

Learning Fair Pareto-Optimal Policies in Multi-Objective Reinforcement Learning

Anonymous authors

Paper under double-blind review

Keywords: Multi-objective reinforcement learning, Deep reinforcement learning, Fair optimization, Welfare functions

Summary

Fairness is important in multi-objective reinforcement learning (MORL), where policies must balance optimality and equity across objectives. While *single-policy* MORL methods can learn fair policies for fixed user preferences using welfare, they fail to generalize for different user preferences. To address this limitation, we propose a novel framework for fairness in *multi-policy* MORL, which learns a set of fair policies. Our theoretical analysis establishes that for concave and piecewise-linear welfare functions, fair policies remain in the convex coverage set (CCS). Additionally, we demonstrate that non-stationary and stochastic policies improve fairness over stationary and deterministic policies. Building on our theoretical analysis, we introduce three scalable methods: an extension of Envelope for fair stationary policies, a non-stationary counterpart using state-augmented accrued rewards, and a novel extension for learning stochastic policies. We validate our methods through extensive experiments across three domains and show that our methods fairer solutions as compared to MORL baselines.

Contribution(s)

1. We introduce a novel framework for fairness in multi-policy MORL, which enables learning a set of fair policies for varying user preferences.
Context: Prior work on fairness in MORL has mainly focused on a single policy for predefined preference weights via some welfare functions. Our framework generalizes fairness across multiple policies, which allow end users to select any policy provided by their preference weights.
2. We provide theoretical analysis demonstrating that for concave, piecewise-linear welfare functions, fair policies remain in the convex coverage set (CCS). Additionally, we establish that non-stationary and stochastic policies can enhance fairness over stationary and deterministic policies, respectively.
Context: Existing work has explored fairness in RL for predefined preference weights but has not theoretically analyzed how non-stationary and stochastic policies can improve fairness for varying preference weights.
3. We propose three scalable methods for learning fair policies in MORL using a single parameterized network: (i) an extension of Envelope (Yang et al., 2019) for learning fair policies, (ii) a non-stationary extension that incorporates state-augmented accrued rewards to adaptively improve fairness, and (iii) a novel stochastic policy learning method that further enhances fairness.
Context: Unlike prior work on MORL, which typically learns Pareto optimal policies, our methods efficiently learn a set of fair policies while maintaining scalability.

Learning Fair Pareto-Optimal Policies in Multi-Objective Reinforcement Learning

Anonymous authors

Paper under double-blind review

Abstract

1 Fairness is an important aspect of decision-making in multi-objective reinforcement
 2 learning (MORL), where policies must ensure both optimality and equity across mul-
 3 tiple, potentially conflicting objectives. While *single-policy* MORL methods can learn
 4 fair policies for fixed user preferences using welfare functions such as the *generalized*
 5 *Gini welfare function* (GGF), they fail to provide the diverse set of policies necessary for
 6 dynamic or unknown user preferences. To address this limitation, we formalize the fair
 7 optimization problem in *multi-policy* MORL, where the goal is to learn a set of Pareto-
 8 optimal policies that ensure fairness across all possible user preferences. Our key tech-
 9 nical contributions are threefold: (1) We show that for concave, piecewise-linear wel-
 10 fare functions (e.g., GGF), fair policies remain in the *convex coverage set* (CCS), which
 11 is an approximated Pareto front for linear scalarization. (2) We demonstrate that non-
 12 stationary policies, augmented with accrued reward histories, and stochastic policies
 13 improve fairness by dynamically adapting to historical inequities. (3) We propose three
 14 novel algorithms, which include integrating GGF with multi-policy multi-objective Q-
 15 Learning (MOQL), state-augmented multi-policy MOQL for learning non-stationary
 16 policies, and its novel extension for learning stochastic policies. To validate the per-
 17 formance of the proposed algorithms, we perform experiments in various domains and
 18 compare our methods against the state-of-the-art MORL baselines. The empirical re-
 19 sults show that our methods learn a set of fair policies that accommodate different user
 20 preferences.

21 1 Introduction

22 Multi-objective reinforcement learning (MORL) is an important topic in the area of reinforcement
 23 learning (RL) that focuses on designing control policies to optimize multiple objectives simultane-
 24 ously. While traditional MORL methods focus on learning Pareto optimal solutions—ensuring no
 25 objective can be improved without sacrificing another—they often neglect fairness, which requires
 26 equitable treatment of all objectives or users in our context. For example, in healthcare, a policy
 27 may aim to maximize overall patient outcomes (optimality) while ensuring equal treatment across
 28 different demographic groups (fairness). A common approach to solving fairness in MORL is to use
 29 *utilitarian* welfare functions, where user utilities are aggregated, typically via weighted sum, into
 30 a scalarized objective. Despite its simplicity, this approach struggles with fairness, as some users’
 31 utilities may be significantly reduced to achieve overall efficiency. An alternative approach is to
 32 employ an *egalitarian* welfare function, which prioritizes the least advantaged user by maximizing
 33 the minimum utility. While this approach improves fairness, it often leads to inefficient solutions
 34 overall, as it optimizes only the lowest utility without ensuring fairness across all objectives.

35 Several works have explored fairness in the *single-policy* RL setting (Weng, 2019; Siddique et al.,
 36 2020; Zimmer et al., 2021; Chen & Hooker, 2021; Do & Usunier, 2022; Fan et al., 2022; Yu et al.,
 37 2023b; Nashed et al., 2023), where only a single policy is learned. For instance, the work in (Weng,

2019) and (Siddique et al., 2020) enforced fairness by utilizing the generalized Gini social welfare function as a scalarized function and assigning appropriate weights to different objectives to ensure their equitable treatment. Extensions have been explored in multi-agent RL (Zimmer et al., 2021; Siddique et al., 2024b) and preferential treatment under known preference weights (Yu et al., 2023b). Recently, fairness has been studied in multi-policy MORL (Cimpeana et al., 2023; Michailidis et al., 2024) where Cimpeana et al. (2023) defined several fairness notions, while (Michailidis et al., 2024) proposed the Lorenz Condition Network (LCN), an extension of the Pareto Conditioned Network (PCN), which trains a policy network in a supervised manner to map states to desired returns. Despite these works, the investigation of fairness in RL still poses some limitations, including (1) learning a *single* fair policy, (2) required knowledge of the welfare function (e.g., scalarized function) with preference weights a priori, and (3) training a conditioning network on specific return targets, limiting their ability to generalize to unseen preferences. Hence, the existing methods operate under fixed/predefined preference weights and cannot be generalized for all possible preferences.

To address these limitations, we propose a novel framework to address fairness in *multi-policy* MORL, rather than the traditional *single-policy* MORL that is the focus of the existing work. Our methods are highly scalable as they leverage a single parameterized network to learn an undominated set of policies, specifically a convex coverage set (CCS), by sampling the entire preference space in MORL. In particular, to address fairness, we apply the welfare function (e.g., GGF) during learning for each sampled preference weight to ensure that each learned policy treats its objectives fairly. We further introduce non-stationary action selection using the state-augmented accrued rewards to enhance fairness by effectively utilizing historical information. We further demonstrate the benefits of learning stochastic policies for fairness. Motivated by hindsight experience replay (Andrychowicz et al., 2017), we incorporate resampling of random preference weights across different preference conditions to improve sample efficiency in MORL, as it is done in (Yang et al., 2019).

The main contributions of this paper are as follows:

1. We introduce a novel framework for fairness in multi-policy MORL, enabling users to select any fair policy based on their specific preferences, thereby enhancing user satisfaction(Section 3.2).
2. We provide theoretical analysis establishing that for concave, piecewise-linear welfare functions (e.g., GGF), fair policies remain in CCS. Additionally, we demonstrate that non-stationary policies can improve fairness by adapting to historical disparities and that stochastic policies further improve fairness over deterministic policies(Section 4).
3. Building on our theoretical insights, we propose three scalable methods for learning fair policies in MORL using a single parameterized network: (i) an extension to Envelope (Yang et al., 2019) for learning fair stationary policies, (ii) a non-stationary counterpart that incorporates state-augmented accrued rewards to adaptively improve fairness over time, and (iii) a novel extension for learning stochastic policies, which further enhances fairness(Section 5).
4. We experimentally validate our methods and demonstrate their effectiveness compared to state-of-the-art MORL and fairness methods across three different domains(Section 6).

2 Related Work

Fairness in machine learning (ML) has become a significant research direction (Dwork et al., 2012; Zafar et al., 2017; Sharifi-Malvajerdi et al., 2019; Singh & Joachims, 2019; Chierichetti et al., 2017; Busa-Fekete et al., 2017; Agarwal et al., 2018; Nabi et al., 2019; Zhang & Liu, 2021). Several studies have addressed fairness in model predictions (Speicher et al., 2018), recommender systems (Leonhardt et al., 2018), classification (Dwork et al., 2012; Zafar et al., 2017; Agarwal et al., 2018; Kim et al., 2019), and ranking (Singh & Joachims, 2019). While much of the literature focuses on the principle of “equal treatment of equals”, other aspects, such as proportionality (Bei et al., 2022) or envy-freeness (Chevalerey et al., 2006) and its multiple variants (e.g., (Beynier et al., 2019; Chakraborty et al., 2021)), have been considered in ML. In contrast, our work is grounded in distributive justice (Rawls, 1971; Brams & Taylor, 1996; Moulin, 2004), with a focus on optimiz-

ing a welfare function for fairness considerations. This principled approach has also been recently advocated in several papers (Heidari et al., 2018; Speicher et al., 2018; Cousins, 2021).

Recently, fairness in RL has gained significant attention with the work by (Jabbari et al., 2017), which ensures fairness in state visitation using scalar rewards. The work of (Jiang & Lu, 2019), proposed FEN a hierarchical decentralized method using a gossip algorithm to ensure fairness across agents involved in a system. Similarly, (Chen et al., 2021) proposed to incorporate fairness into actor-critic RL algorithms, optimizing general fairness utility functions for real-world network optimization problems. Considering the multi-objective nature of many RL problems, the study of fairness in multi-objective reinforcement learning (MORL) has been widely studied. In particular, (Siddique et al., 2020) proposed multiple adaptations to deep RL algorithms that optimize the *generalized Gini social welfare*. (Zimmer et al., 2021; Siddique et al., 2024a) extended this work to the decentralized cooperative multi-agent setting. (Fan et al., 2022) proposed to optimize the Nash welfare function using scalarized expected return criterion. (Do & Usunier, 2022) proposed a method for generalized Gini welfare function optimization in rankings. (Yu et al., 2023b; Qian et al., 2025) proposed methods that learn a fair policy providing preferential treatment to some users while ensuring equal treatment of all others under the assumption that these preferential weights are known in advance. (Siddique et al., 2023) proposed FPbRL, a fairness-enhanced method in preference-based RL to learn fair policies in the absence of true rewards. Recently, fairness has been considered in multi-policy MORL with (Michailidis et al., 2024) propose learning Lorenz Condition networks, which ensures fairness through Lorenz domination and adds an extra parameter λ , however, we use the welfare function to learn a set of fair optimal policies.

Despite the significant successes achieved in the field of deep RL and MORL, existing methods heavily rely on scalarization functions to learn a *single policy* with fixed preference weights. However, such single-policy methods do not work when preferences are unknown or user-specific solutions are required. To address this limitation, several works have been proposed to accommodate user-specific preferences, including but not limited to those proposed by (Barrett & Narayanan, 2008; Van Moffaert et al., 2013; Moffaert & Nowé, 2014; Yang et al., 2019; Alegre et al., 2023; Reymond et al., 2022). Notably, these methods aim to learn a set of policies that approximate the Pareto frontier of optimal solutions. For instance, (Barrett & Narayanan, 2008) and (Moffaert & Nowé, 2014) proposed methods to compute policies on the Pareto front’s convex hull, while (Yang et al., 2019) introduced envelope Q-learning, learning policies from the convex coverage set (CCS). These approaches, however, do not address fairness, which is the focus of this paper.

3 Preliminaries

3.1 Multi-Objective Markov Decision Process

A multi-objective Markov Decision Process (MOMDP) extends the classical Markov Decision Process (MDP) framework to scenarios where an agent must optimize multiple objectives simultaneously. An MDP (Puterman, 1994) is a mathematical model commonly used for sequential decision-making problems. Formally, an MDP is defined by a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$, where \mathcal{S} is the set of states, \mathcal{A} is the set of actions available to the agent, $\mathcal{P}_{a,s,s'} \in [0, 1]$ is the probability of transition from state s to state s' after taking action a , i.e., $\mathcal{P}(s'|s, a) = \mathcal{P}[S_{t+1} = s'|S_t = s, A_t = a]$, $r(s, a) : s \times a \mapsto r$ is the immediate reward obtained by taking action a at state s , and $\gamma \in [0, 1]$ is the discount factor. An MOMDP can be represented by a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathbf{r}, \gamma, \Omega, f_\Omega)$, in which the definitions of $\mathcal{S}, \mathcal{A}, \mathcal{P}$, and γ are the same as in MDP except that the reward \mathbf{r} is now a vector, with each component corresponding to an objective that the agent seeks to optimize. Here, the additional Ω represents the entire space of preferences, and f_Ω is the preference function which takes a linear form, producing a single utility $f_\omega(\mathbf{r}) = \omega^T \mathbf{r}(s, a)$, where ω is a vector representing the preference weights for different objectives. In MOMDPs, the objectives may be conflicting, and hence it is often difficult to optimize all objectives simultaneously.

The goal of an agent in an MOMDP is to either learn a single policy that balances multiple objectives or a set of policies that optimize different trade-offs among objectives. These approaches are referred to as *single-policy* MORL and *multi-policy* MORL, respectively. A policy π is a strategy that maps states to actions, which can be deterministic (i.e., $\forall s, \pi(s) \in \mathcal{A}$) or stochastic (i.e., $\forall s, a, \pi(a|s)$ denotes the probability of selecting a in s). In MOMDPs, policies are typically *stationary* or *Markovian*, meaning that action selection probabilities depend solely on the current state, irrespective of past states and actions. Conversely, a non-stationary policy $\pi(a|\tau, s)$, also known as an adaptive policy, may depend on the agent’s history τ . Standard definitions in MDPs, such as the return $G(\tau)$ and the value functions V or Q , extend naturally to MOMDPs, albeit represented as vectors and matrices respectively. The vector return in an MOMDP is expressed as $\mathbf{G}(\tau) = \sum_{t=1}^{\infty} \gamma^{t-1} \mathbf{r}_t$, where τ is a trajectory comprising a sequence of states, actions, and rewards following the policy, and \mathbf{r}_t is a vector reward obtained at time step t . The state value function of a policy π in an MOMDP is defined as $\mathbf{V}^\pi(s) = [V_i^\pi(s)] = \mathbb{E}_{\tau \sim \pi} [\sum_{t=1}^{\infty} \gamma^{t-1} \mathbf{r}_t \mid S_0 = s]$, where all operations (addition, product) are applied component-wise.

In MOMDPs, value functions do not offer a complete ordering over the policy space. This means it is possible to encounter scenarios wherein, e.g., $V_i^\pi(s) > V_i^{\pi'}(s)$ for objective i , while $V_j^\pi(s) < V_j^{\pi'}(s)$ for objective j . Hence, value functions in MOMDPs induce only a partial ordering within the policy space, necessitating additional information into objective prioritization for policy ordering.

Envelope Multi-Objective Q-Learning. The Envelope algorithm (Yang et al., 2019) learns a convex coverage set (CCS) by sampling preference weights $\omega \in \Omega$ and optimizing linearly scalarized Q-values: $Q(s, a, \omega) = \omega^T \mathbf{Q}(s, a)$, where $\mathbf{Q}(s, a) \in \mathbb{R}^N$ is the vector of Q-values for N objectives. The Bellman optimality equation for Envelope algorithm is: $\mathbf{Q}^*(s, a, \omega) = \mathbf{r}(s, a) + \gamma \max_{a'} \omega^T \mathbf{Q}^*(s', a')$. A single neural network parameterizes $\mathbf{Q}(s, a, \omega)$ by concatenating ω to the state s , enabling efficient learning across all preferences. Despite its scalability, Envelope lacks explicit fairness guarantees, as linear scalarization may prioritize dominant objectives.

3.2 Fairness Formulation

In MORL, fairness, rooted in distributive justice (Moulin, 2004), is crucial for ensuring equitable distribution of rewards. Prior studies in fair optimization within MORL have primarily focused on learning a *single-policy*, commonly referred to as an average policy (Siddique et al., 2020; Fan et al., 2022; Yu et al., 2023a; Siddique et al., 2023). In this paper, we adopt a more inclusive view of fairness, including *efficiency*, *equity*, and *impartiality* to generate fair optimal solutions for user-specific preferences. For discussion on fairness and welfare function, please refer to the Appendix.

Definition 3.1. *Efficiency states that among two solutions, if one solution is (weakly or strictly) preferred by all users, then it should be preferred to the other one, e.g., $\mathbf{V} \succ \mathbf{V}' \Rightarrow \phi(\mathbf{V}) > \phi(\mathbf{V}')$, where $\phi(\mathbf{V})$ is the scalar utility function by using the ϕ that specifies the value of a solution.*

The efficiency property specifies that given all else equal, one prefers to increase a user’s utility. In the MORL setting, the efficiency property simply means Pareto dominance. More specifically, a solution is considered efficient if it is not dominated by any other solution for all objectives.

Definition 3.2. *For a given pair of solutions $\mathbf{V}, \mathbf{V}' \in \mathbb{R}^N$, \mathbf{V} weakly Pareto-dominates \mathbf{V}' if $\forall i, V_i \geq V'_i, \forall i \in \{1, \dots, N\}$, where N is the total number of objectives. Besides, \mathbf{V} Pareto-dominates \mathbf{V}' if $V_i \geq V'_i, \forall i$ and $\exists j, V_j > V'_j$. For brevity, we denote Pareto dominance as \geq for the weak form and $>$ for the strict form.*

Essentially, a solution \mathbf{V} (weakly) Pareto-dominates another solution \mathbf{V}' if the former’s value $\phi(\mathbf{V})$ (weakly) Pareto-dominates that of the latter $\phi(\mathbf{V}')$. A solution \mathbf{V}^* is said to be *Pareto-optimal* if no other solution \mathbf{V} Pareto-dominates it. *Pareto front* (\mathcal{F}) is defined as the set of Pareto-optimal solutions, which may consist of infinitely many solutions, especially when policies can be stochastic. A typical way to approximate (\mathcal{F}) is to compute the convex coverage set (CCS), defined below.

Definition 3.3. A solution in CCS has a maximal scalarized value in a weighted sense if there exists a weight vector $\omega \in \Omega$ such that the scalarized utility $\omega^T \mathbf{V}$ is weakly preferred to the scalarized utility $\omega^T \mathbf{V}'$ for all other solutions \mathbf{V}' in the Pareto front. Formally speaking, $\mathbf{V} \in \text{CCS} \iff \exists \omega \in \Omega$ s.t. $\omega^T \mathbf{V} \geq \omega^T \mathbf{V}', \forall \mathbf{V}' \in \mathcal{F}$.

Next, we discuss the significance of the *equity* property, a stronger property than efficiency and often associated with distributive justice, as it refers to the fair distribution of resources or opportunities. This property ensures that a fair solution follows the *Pigou-Dalton principle* (Moulin, 2004), which states the transferring of rewards from more advantaged users to less advantaged users.

Definition 3.4. A solution satisfies the Pigou-Dalton principle if for all \mathbf{V}, \mathbf{V}' equal except for $V_i = V'_i + \delta$ and $V_j = V'_j - \delta$ where $V'_i - V'_j > \delta > 0$, $\phi(\mathbf{V}) > \phi(\mathbf{V}')$.

Finally, the *impartiality* property, which is rooted in the principle of “equal treatment of equals” states that individuals sharing similar characteristics should be treated similarly.

Definition 3.5. In a system, individuals with similar characteristics should be treated similarly, i.e., the solution should be independent of the order of its arguments $\phi(\mathbf{V}) = \phi(\mathbf{V}_\sigma)$, where σ is a permutation and \mathbf{V}_σ is the vector obtained from vector \mathbf{V} permuted by σ .

To ensure fairness that satisfies the above three properties, we use a well-known generalized Gini welfare function (GGF) (Weymark, 1981), which can be defined as:

$$\phi_{\text{GGF}}(\mathbf{u}) = \sum_{i \in N} \omega_i u_i^\uparrow, \quad (1)$$

$\mathbf{u} \in \mathbb{R}^N$ represents the utility vector of a size N for N objectives, $\omega \in \mathbb{R}^N$ is a fixed weight vector with positive components that strictly decrease (i.e., $\omega_1 > \dots > \omega_N$) with $\sum_i \omega_i = 1$, and \mathbf{u}^\uparrow denotes the vector by sorting the components of \mathbf{u} in an increasing order (i.e., $u_1^\uparrow \leq \dots \leq u_N^\uparrow$). GGF satisfies the aforementioned three fairness properties. As the weights are positive, it is monotonic with respect to Pareto dominance, thus satisfying the efficiency property. Since the utility vector is reordered, it is also symmetric and therefore satisfies the impartiality property. Furthermore, the positive and decreasing weights ensure that GGF is Schur-concave, i.e., monotonic with respect to Pigou-Dalton transfers, therefore satisfies the impartiality property.

GGF has been studied and used in MORL extensively (Siddique et al., 2020; Mandal & Gan, 2022; Yu et al., 2023a; Qian et al., 2025), however, all of these works used it for single-policy setting. We are the first ones to use it in a multi-policy MORL setting. In multi-policy MORL, the usual approach is to find all Pareto non-dominated solutions (Mukai et al., 2012; Van Moffaert & Nowé, 2014). This approach may work for small problems, however, for large-scale problems, the Pareto non-dominated solutions grow exponentially. A better way to achieve scalable and multiple solutions to approximate the Pareto front is possibly to arrive at the solutions that form the convex envelope and thus form a convex coverage set.

4 Fairness in MORL

Since we are in a multi-policy MORL setting, where an agent learns a set of Pareto optimal policies, fairness becomes more important as different stakeholders may have different preferences and during inference, any solution can be used from the Pareto non-dominated solutions given the stakeholder preferences. We formalize this sophisticated multi-policy fair optimization problem as:

$$\forall \omega \in \Omega, \quad \max_{\pi \in \Pi} \phi_{\text{GGF}}(\mathbf{J}(\pi)), \quad (2)$$

where Ω is the set of valid preference weights sorted in descending order, $\mathbf{J}(\pi) = \mathbb{E}_\pi[\sum_{t=0}^{\infty} \gamma^t \mathbf{r}_t]$ is the expected discounted return, and $\phi_{\text{GGF}}(\mathbf{J}) = \sum_{i=1}^N \omega_i J_{(i)}$ with $J_{(1)} \leq \dots \leq J_{(n)}$. The concavity of GGF makes problem (2) as convex optimization problem, enabling efficient solutions within the

CCS. Below, we establish three foundational results, which show that it is always feasible to obtain optimal solutions in the CCS corresponding to GGF fair optimization. Next, we demonstrate that a non-stationary policy based on accrued rewards is beneficial in yielding improved fairness when compared with its stationary counterpart. Here, a policy yields improved fairness or is fairer if a higher welfare score, defined in (1), is achieved. Lastly, we show that a stochastic policy may yield fairer solutions than a deterministic one.

Sufficiency of Optimal Solutions in the CCS. The first question relates to the learning of fair policies in a multi-policy MORL setting is which subset of policies may be optimal among the set of all (possibly non-stationary) policies. Indeed, for linear scalarization function, CCS contains the set of Pareto front solutions. Below, we formally state it:

Lemma 4.1. *For any MOMDP with linear preferences over objectives, the CCS contains an optimal policy for any linear combination of the objectives.*

While GGF introduces non-linear fairness objectives, its piecewise linearity and concavity allow representation as a maximum of linear functions, which ensures that solutions lie within the CCS. The following proposition establishes the sufficiency of the CCS in representing optimal policies for ϕ_{GGF} preference weights.

Proposition 4.1. *For any $s \in \mathcal{S}$ in an MOMDP and a piecewise-linear concave welfare function ϕ_{GGF} (e.g., GGF) that can be represented as, $\phi_{GGF}(\mathbf{V}^\pi(s)) = \min_{\sigma \in \mathbb{S}_N} \{\omega_\sigma^\top \mathbf{V}^\pi(s)\}$, there exists a policy $\pi^* \in \text{CCS}$ such that $\phi_{GGF}(\mathbf{V}^{\pi^*}(s)) \geq \phi_{GGF}(\mathbf{V}^\pi(s))$, $\forall \pi \in \Pi$.*

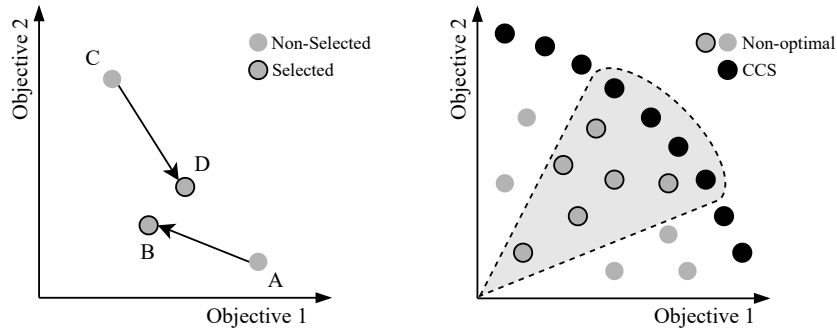


Figure 1: Examples of 2-objective MOMDP where GGF leads to fairer outcomes.

Example 4.1 To illustrate how the GGF function ensures fairness in MORL, consider a two-objective MOMDP with objective values $\mathbf{V}_1 = (3, 1)$ and $\mathbf{V}_2 = (2, 3)$ and weights $(1, 2)$. For \mathbf{V}_1 , two weighted combinations are possible: **A**) $(3, 1) \cdot (2, 1) = (6, 1)$ with scalar sum $6 + 1 = 7$, **B**) $(3, 1) \cdot (1, 2) = (3, 2)$ with scalar sum $3 + 2 = 5$. Since the GGF is defined as $\phi_{GGF}(\mathbf{V}^\pi(s)) = \min_{\sigma \in \mathbb{S}_N} \{\omega_\sigma^\top \mathbf{V}^\pi(s)\}$, it selects the lower scalar value, preferring point B over A (see left figure of Figure 1). Similarly, for \mathbf{V}_2 : **C**) $(2, 3) \cdot (1, 2) = (2, 6)$ with scalar sum $2 + 6 = 8$, **D**) $(2, 3) \cdot (2, 1) = (4, 3)$ with scalar sum $4 + 3 = 7$. Here, point D is preferred over C. This mechanism directs the solutions toward the fairer region (grey dotted area in the right figure of Figure 1), demonstrating that maximizing the GGF leads to fair Pareto-optimal solutions.

Fairness of Non-Stationary Policies. In fair MORL, learning non-stationary policies can be particularly beneficial, as they leverage historical information to make more informed decisions and adapt over time.

Proposition 4.2. *Let the reward \mathbf{r} be nonnegative, and Π_S and Π_{NS} be the sets of stationary and non-stationary policies, respectively. For any $s \in \mathcal{S}$ in an MOMDP and a given ϕ_{GGF} , there exists a non-stationary policy $\pi_{NS} \in \Pi_{NS}$ that achieves a higher welfare score than any stationary policy $\pi_S \in \Pi_S$, i.e., $\exists \pi_{NS} \in \Pi_{NS} : \phi_{GGF}(\mathbf{V}^{\pi_{NS}}(s)) \geq \max_{\pi_S \in \Pi_S} \phi_{GGF}(\mathbf{V}^{\pi_S}(s))$.*

Example 4.2 To illustrate the value of learning a non-stationary policy, consider a 2-objective MOMDP, shown in Fig. 2. At timestep $t > 0$, the agent has accrued a vector reward $\mathbf{r}_{acc} = (10, 0)$ for two objectives. The preference weights, encapsulated within the welfare function ϕ , denote decreasing weights, such as $(0.8, 0.2)$. With two potential actions, each leading to a final state, action a_1 yields a reward of $(0, 10)$, while action a_2 yields $(5, 5)$. Since s_t is the absorbing state, we can set discount factor $\gamma = 1$. Under the given welfare function ϕ defined in 1, executing a_1 yields a welfare score of 2, whereas executing a_2 yields a score of 5 if only future rewards are considered. However, considering historical data, i.e., \mathbf{r}_{acc} , a_1 yields a higher accrued episodic return of $(10, 10)$ and a welfare score of 10. Similarly, a_2 yields $(15, 5)$ and 7 episodic return and welfare scores, respectively. Note that action a_1 is a fairer choice in this case since it balances the two objectives, unlike action a_2 , which fails to achieve a more equitable outcome. Hence, employing historical data, namely, accrued rewards in this case, is critical to enable fair policy learning.

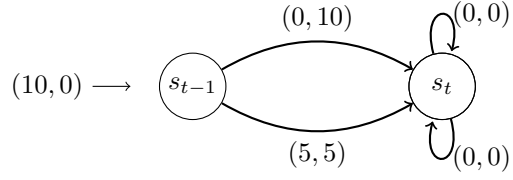


Figure 2: Example of MOMDP where actions lead to different rewards.

Optimality of Stochastic Policies for Fairness Unlike single-objective RL, in MORL, a deterministic policy may not be optimal. A fairer solution can often be achieved through randomization.

Proposition 4.3. Let Π_{ST} be the set of stochastic policies and Π_D be the set of deterministic policies. For an MOMDP \mathcal{M} and a concave welfare function such as ϕ_{GGF} , there exists a stochastic policy $\pi_{ST} \in \Pi_{ST}$ such that $\phi_{GGF}(\mathbf{V}^{\pi_{ST}}) \geq \max_{\pi_D \in \Pi_D} \phi_{GGF}(\mathbf{V}^{\pi_D})$.

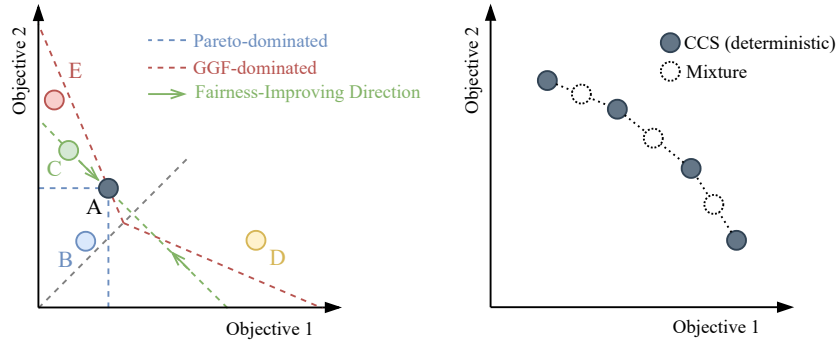


Figure 3: Left Figure: Point A is always preferred to B due to Pareto dominance; A is always preferred to C due to the Pigou-Dalton transfer principle (fairer solution); depending on the weights of GGF, Points D and E can be dominated or non-dominated by A (w.r.t. GGF); with weights $(0.3, 0.7)$, A is preferred to E but not to D. Right Figure: Black points refer to deterministic policy that in CCS and stochastic policy can be obtained with the mixture of deterministic policies in the CCS, shown in dotted point. Demonstrate stochastic policy can achieve fairer solution which deterministic policy cannot.

The proofs of the above lemma and propositions are provided in Section 8. Left figure of Figure 3 illustrates GGF on a two-objective optimization task. The optimality of stochastic policies implies that restricting the search for fair solutions to deterministic policies is insufficient. Stochastic policies offer a broader range of solutions and may better capture the trade-offs among multiple objectives, enhancing the overall fairness of the policy, shown in the right figure of Figure 3.

5 Proposed Algorithms

In this section, we introduce three novel algorithms that incorporate fairness into MORL based on the technical analysis in the previous section. These algorithms optimize the GGF welfare function defined in (1) to ensure fairness across N fixed users with varying preferences. Our proposed methods are scalable and sample-efficient as they utilize a single parameterized network to estimate Q-values for all objectives while maintaining a diverse set of Pareto-optimal policies. We present three distinct algorithms: Fair Multi-Objective Deep Q-Learning (F-MDQ), its extension with non-stationary policies (FN-MDQ), and a novel extension incorporating stochastic policies (FNS-MDQ). This progression from stationary to non-stationary to stochastic and non-stationary policies demonstrates our systematic approach to enhancing fairness in MORL algorithms, with each method building upon and improving the previous one.

F-MDQ. F-MDQ builds on the Envelope algorithm (Yang et al., 2019) by replacing the linear scalarization function with the GGF welfare function ϕ . This ensures fairness while learning policies across all preferences $\omega \in \Omega$. The Bellman optimality equation for F-MDQ is given by:

$$Q^*(s, a, \omega) = \mathbb{E}[r(s, a) + \gamma Q^*(s', \sup_{a' \in \mathcal{A}} \phi_{\text{GGF}}(r(s, a) + Q^*(s', a', \omega), \omega) \mid s, a)],$$

where $Q^\pi(s, a, \omega)$ represents the expected return vector for policy π , conditioned on preference ω . As the MO Q-function is parameterized, it can be learned by minimizing the loss function $\mathcal{L} = \mathbb{E}_{(s, a, r, s', \omega) \sim \mathcal{D}} [\|y - Q(s, a, \omega)\|_2^2]$, where the expectation is taken over experiences sampled from the replay buffer \mathcal{D} . Given that the loss function includes an expectation over ω , the preference weights are sampled randomly and are decoupled from the transitions, allowing increased sample efficiency through a resampling scheme similar to Hindsight Experience Replay (HER) (Andrychowicz et al., 2017). The target y for F-MDQ is computed as $y = r(s, a) + \gamma Q'(s', \sup_{a' \in \mathcal{A}} \phi_{\text{GGF}}(r(s, a) + \gamma Q(s', a', \omega), \omega))$, where Q' represents the target multi-objective Q-function, and the supremum is applied over the GGF welfare function ϕ_{GGF} instead of a linear weighted sum. This ensures that actions are selected based on higher welfare scores rather than simply maximizing Q-values.

FN-MDQ. FN-MDQ extends F-MDQ by incorporating accrued rewards into the state to learn non-stationary policies, as discussed in Proposition 8.2. It augments the observed state with accrued rewards, allowing the agent to balance reward distribution across users more effectively (as demonstrated in Example 2). The augmented state is defined as $\bar{s}_t = (s_t, r_{\text{acc}})$, where $r_{\text{acc}} = \sum_{i=1}^{t-1} \gamma^{i-1} r_i$ is the discounted total reward received in the current trajectory. The regression target for FN-MDQ is then given by $r(s_t, a_t) + \gamma Q'(\bar{s}_{t+1}, \sup_{a' \in \mathcal{A}} \phi_{\text{GGF}}(Q(\bar{s}_{t+1}, a', \omega)), \omega)$. Here, the immediate reward $r(s_t, a_t)$ is excluded from the optimal action computation since this signal is already included in the augmented state as part of the discounted total reward. This extension enables the agent to identify and prioritize users who have received insufficient rewards within an episode.

FNS-MDQ. Given that stochastic policies can outperform deterministic ones (as established in Proposition 8.3), the performance of FN-MDQ can be enhanced by incorporating stochastic policies. We now explain how stochastic policies can be integrated into the FN-MDQ algorithm.

Under the stochastic policies, the target Q-value is adjusted to account for the expected Q-values, which reformulates the update as $r(s_t, a_t) + \gamma Q'(\bar{s}_{t+1}, \sum_{a' \in \mathcal{A}} \phi_{\text{GGF}}(\pi(a' \mid \bar{s}_{t+1}) Q(\bar{s}_{t+1}, a', \omega)), \omega)$, where $\pi(a' \mid \bar{s}_{t+1})$ is the probability of taking action a' given the augmented state \bar{s}_{t+1} . This reformulation considers the distribution of possible actions rather than selecting a single best deterministic action, aligning with our theoretical insights.

Unlike F-MDQ and FN-MDQ, which rely on deterministic action selection, FNS-MDQ samples actions from a probability distribution over Q-values. This stochastic action selection improves fairness by enabling more balanced policy exploration and reducing biases that arise from always selecting the highest Q-value action. Note that, during the training phase, all algorithms employ

an ϵ -greedy policy during training, however, FNS-MDQ differs in its action-selection strategy by using the best learned stochastic policy rather than a deterministic greedy approach. This increased flexibility and randomness can lead to more equitable solutions.

6 Experiments

To evaluate the proposed methods, we conduct experiments across three domains—each characterized by varying levels of complexity in terms of the number of objectives. These domains, ranging from low to high in terms of the number of objectives, include species conservation, resource gathering, and multi-product web advertising. Each environment presents unique challenges where fairness plays a critical role. We first briefly describe each environment (details are available in the Appendix B) and then present our experimental results.

6.1 Environments

Our first domain is a species conservation (SC) environment, which addresses a critical ecological challenge: balancing the populations of two highly interacting endangered species, the sea otter and the northern abalone. Both species are at risk of extinction, requiring sophisticated management strategies to ensure their survival. We adopt the model proposed by (Chadès et al., 2012), which simulates the predation relationship between the species, where sea otters prey on abalones. This dynamic presents a unique preservation challenge, as the survival of one species could potentially drive the other to extinction if not properly managed. The state space is composed of the current population sizes of sea otters and northern abalones. The action space includes introducing sea otters, enforcing anti-poaching measures, controlling sea otter populations, implementing a combination of half-antipoaching and half-controlled sea otters, or taking no action. Each action has significant ecological implications. For instance, introducing sea otters may help balance the abalone population, but if mismanaged, could lead to abalone extinction. The reward function is defined by the population densities of both species, i.e., $N = 2$. Fairness in this context is interpreted as achieving a balanced distribution of species densities to ensure their preservation.

Our second environment is a resource-gathering (RG) problem, which is a 5×5 grid world that contains three types of resources: gold, gems, and stones. These resources are randomly positioned on the grid and regenerate randomly upon consumption. The main challenge here is to collect these resources, where each resource has a different value: gold and gems are valued at 1, while stones have a lower value of 0.4. This creates an intentionally uneven resource distribution, with two stones, one gold, and one gem. In this environment, the state is defined by the agent’s current location on the grid and the cumulative count of each resource collected during its trajectory. The agent can take four actions: up, down, left, and right. The reward function is defined as a vector representing the rewards collected for each type of resource. In this environment, fairness is defined as the equitable collection of resources, despite their differing values. Note that, this problem is particularly important for validating whether the proposed methods can achieve fairer solutions while still reaching Pareto optimal solutions.

Our third domain is a multi-product web advertising (MWP) problem that involves an online store offering $N = 7$ distinct products. Here, the agent decides which advertisement to display: a product-specific advertisement for one of the products $i \in [0, \dots, N - 1]$, or a general advertisement that is not tailored to any specific product. In this environment, the state space includes the number of products available in the store, as well as the number of visits, purchases, and exits. The action space is $N + 1$, where actions 0 through $N - 1$ correspond to displaying advertisements for specific products, and action N involves showing a general advertisement. This additional action adds complexity, requiring the agent to decide the optimal moment to transition between states. The reward function is designed so that the agent receives a reward of 1 in the i^{th} dimension of the reward vector if a product of the type i is sold after displaying its advertisement. In this environment, fairness is defined as balancing the frequency of advertisements shown for each product, ensuring no single

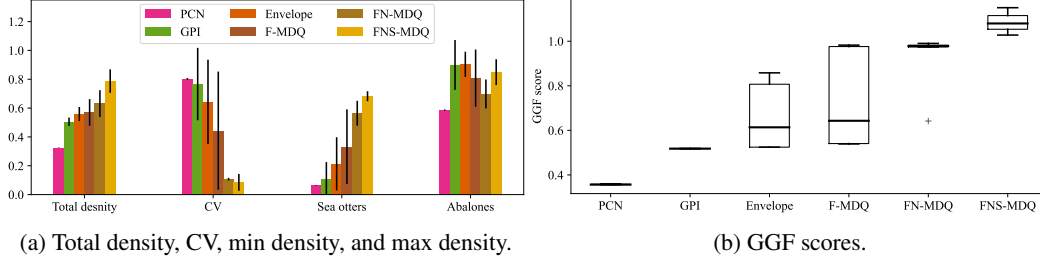


Figure 4: Performances of multi-policy MORL baselines and our methods in species conservation.

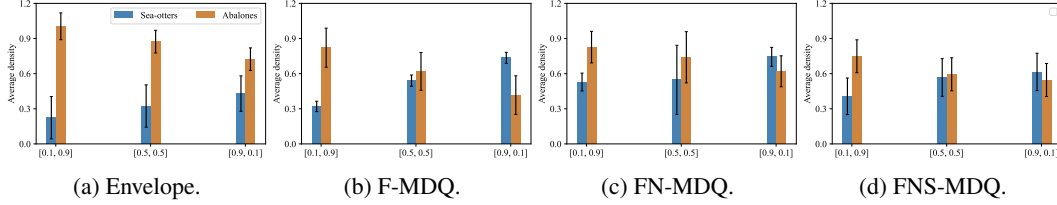


Figure 5: Individual densities of Envelope, and our proposed methods during testing with unseen preferences in species conservation.

380 product is overly prioritized. The challenge lies in increasing overall rewards while maintaining a
 381 fair distribution of advertisement exposure across all products.

382 6.2 Baselines

383 We compare our proposed methods against several multi-policy MORL baselines. Generalized Pol-
 384 icy Improvement Linear Support (GPI-LS) (Alegre et al., 2023) employs GPI (Barreto et al., 2017)
 385 to combine policies within its learned Convex Coverage Set (CCS) and prioritize the weight vectors
 386 on which agents should train at each moment. The Envelope algorithm (Yang et al., 2019) uses a
 387 single neural network conditioned on a weight vector to approximate the CCS. Pareto Conditioned
 388 Networks (PCN) (Reymond et al., 2022) utilizes a neural network conditioned on a desired return
 389 per objective and is trained via supervised learning to predict actions that yield the desired return.
 390 Hyperparameters for each method were optimized, and experiments were run for five different seeds,
 391 with average results reported. Further details on experimental configurations and hyperparameters
 392 are provided in Appendix C.

393 6.3 Results

394 In this section, we present the experimental results across the three environments presented above.
 395 The primary objective of these experiments is to assess the effectiveness of our proposed methods
 396 by addressing the following key research questions: **(A)** How effective are our methods in learning
 397 fairer solutions compared to multi-policy MORL baselines? **(B)** Can our methods generate fair solu-
 398 tions across different preference settings during inference? **(C)** What is the impact of our approach
 399 on the diversity and quality of non-dominated solutions that satisfy fairness criteria? **(D)** Does the
 400 incorporation of stochastic policies in MO Q-learning based algorithms contribute to improved fair-
 401 ness or overall performance?

402 **Question (A)** To evaluate how effective our methods are in learning fair solutions, we conducted
 403 experiments in the SC, RG, and MWP domains, as shown in Figures 4a, 6a and 7a. We compare our
 404 proposed methods (F-MDQ, FN-MDQ, and FMS-MDQ) with multi-policy MORL baselines such as
 405 PCN, GPI, and Envelope during the training phase. We choose these baselines as they are the current
 406 state-of-the-art MORL baselines. The Key evaluation metrics used include total rewards, Coefficient

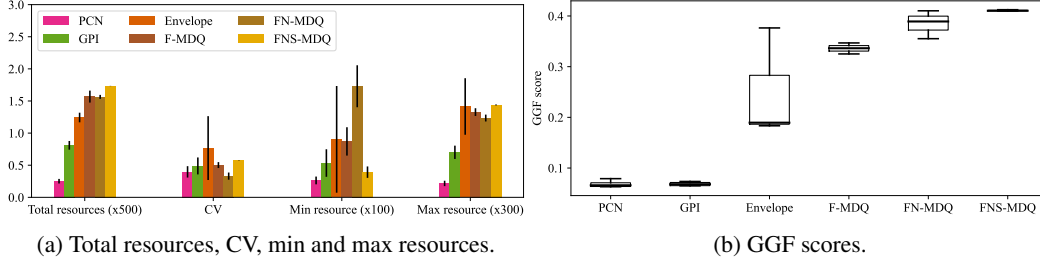


Figure 6: Performances of multi-policy MORL baselines and our methods in resource gathering.

of Variation (CV) indicating the variations in different objectives’ utilities, and the minimum and maximum objective utilities. Moreover, GGF welfare scores were computed to quantify fairness. As we are in a multi-policy MORL, an agent learns a set of Pareto optimal policies during learning. To show the results, we computed these metrics over the last 50 trajectories for all the Pareto optimal policies and reported their normalized scores. Note that, during the last 50 trajectories, all the agents are converged so it ensures a fair comparison for multi-policy MORL methods.

As shown in Figure 4a, PCN performs the worst. GPI outperforms PCN, likely due to its TD3-based (Fujimoto et al., 2018) architecture and efficient prioritization scheme in learning the Pareto front \mathcal{F} . The Envelope algorithm performs better than PCN and GPI as it achieves higher total density and, interestingly, lower CV. However, our proposed algorithms outperform all other methods by achieving the lowest CV and highest welfare scores Figure 4b, with FN-MDQ outperforming F-MDQ, underscoring the value of non-stationary policies. Furthermore, FNS-MDQ outperforms both F-MDQ and FN-MDQ as it maximizes the minimum objective utility and demonstrates better fairness through optimizing the welfare function ϕ_{GGF} . Similar results are observed in RG Figure 6a, where PCN performs the worst as it collects the least resources, likely due to its limitations in deterministic environments (Reymond et al., 2022). Although GPI performs better than PCN, both exhibit low CV alongside poor overall performance and GGF welfare utility Figure 6b. The Envelope algorithm achieves better performance in terms of rewards but suffers from the highest CV and lower GGF utility scores. In contrast, our proposed methods attain a lower CV compared to all baselines, and they achieve the highest GGF scores, highlighting their effectiveness in identifying fair policies through welfare function optimization. Interestingly, FNS-MDQ exhibits a higher CV due to its higher maximum objective and the total resources collected. Nevertheless, it also achieves the highest welfare scores. Consistent with our previous results, our proposed methods in MVP environment Figure 7a achieve the highest welfare scores, indicating their capacity to ensure an equitable distribution of rewards across all objectives. Moreover, they maintain the lowest CV, highlighting their robustness in learning fair policies, even in highly stochastic environments with a higher number of objectives. Once again, PCN, and GPI perform the worst, further underscoring the efficacy of our methods in this context.

Question (B) To check whether our methods can generate fair solutions across different preference settings, we evaluated our algorithms with unseen preferences during testing in the SC environment. As shown in Figure 5, which presents the individual species densities (sea otters and abalones) for preference configurations (0.1, 0.9), (0.5, 0.5), (0.9, 0.1), the Envelope algorithm fails to produce fair solutions, suggesting its limitation in generating fair optimal policies across varying preferences. In contrast, F-MDQ generates more balanced solutions, while FN-MDQ and FNS-MDQ achieve even fairer outcomes, further validating our earlier findings.

Question (C) The results discussed in previous questions suggest that our methods can generate a range of Pareto optimal non-dominated solutions across varied preference configurations, which indicates better coverage of the objective space, thus leading to improved performance across multiple objectives. For quality, our proposed algorithms consistently achieve the lowest CV and highest GGF welfare scores across SC, RG, and MVP domains, indicating that our solutions exhibit more

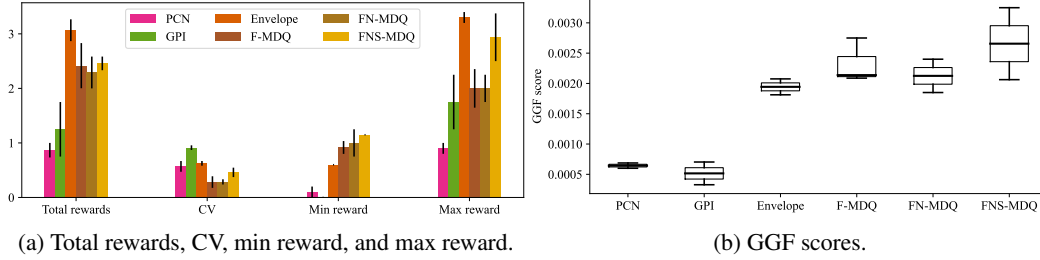


Figure 7: Performances of multi-policy MORL baselines and our proposed methods in the MPW.

equitable distribution of objective utilities while maintaining Pareto optimality compared to baseline methods (PCN, GPI, and Envelope). These outcomes align with our theoretical justifications (see Section 4).

Question (D) Finally, to assess the impact of incorporating stochastic policies in MO Q-learning algorithms, we refer to the results in Figures 4a, 6a and 7a, where stochastic policies consistently improve both efficiency and fairness. Moreover, as shown in Figures 4b, 6b and 7b incorporating stochastic policies also enhances MORL metrics, validating the contribution of stochasticity to both fairness and overall performance.

7 Conclusions and Limitations

In this paper, we presented a novel approach to addressing fairness in the context of multi-policy MORL. Our proposed methods leverage a single parameterized network to learn optimized policies across the entire space of possible preferences. Both theoretical and empirical analyses demonstrate that learning a non-stationary policy significantly improves fairness. Additionally, we highlighted the importance of stochastic policies in achieving fair outcomes. Experimental evaluations in three domains validated the effectiveness of our approach in yielding more equitable policies compared to state-of-the-art MORL and fair baselines.

Our approach also has some limitations. First, it is limited to MOMDPs with discrete action spaces. Second, it assumes that preference weights are linear to learn the CCS, which may not capture the concave regions of the Pareto front. Third, the current formulation is focused on individual fairness. Given that optimizing a welfare function is a broad framework applicable to various real-world MORL problems involving general utilities, an important direction for future research is to extend this approach to accommodate more sophisticated objective functions, particularly those related to group-level fairness, safety, and risk sensitivity.

References

- Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna Wallach. A reductions approach to fair classification. In *ICML*, 2018.
- Lucas N Alegre, Ana LC Bazzan, Diederik M Roijers, Ann Nowé, and Bruno C da Silva. Sample-efficient multi-objective learning via generalized policy improvement prioritization. *arXiv preprint arXiv:2301.07784*, 2023.
- Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. *NeurIPS*, 30, 2017.
- André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt, and David Silver. Successor features for transfer in reinforcement learning. *NeurIPS*, 30, 2017.

- 481 Leon Barrett and Srini Narayanan. Learning all optimal policies with multiple criteria. In *ICML*,
482 2008.
- 483 X. Bei, S. Liu, C.K. Poon, and H. Wang. Candidate selections with proportional fairness constraints.
484 In *AAMAS*, 2022.
- 485 Aurélie Beynier, Yann Chevaleyre, Laurent Gourvès, Ararat Harutyunyan, Julien Lesca, Nicolas
486 Maudet, and Anaëlle Wilczynski. Local envy-freeness in house allocation problems. *AAMAS*,
487 2019.
- 488 Steven J. Brams and Alan D. Taylor. *Fair Division: From Cake-Cutting to Dispute Resolution*.
489 Cambridge University Press, March 1996.
- 490 Róbert Busa-Fekete, Balázs Szörényi, Paul Weng, and Shie Mannor. Multi-objective bandits: Opti-
491 mizing the generalized gini index. In *ICML*, pp. 625–634, 2017.
- 492 Iadine Chadès, Janelle MR Curtis, and Tara G Martin. Setting realistic recovery targets for two
493 interacting endangered species, sea otter and northern abalone. *Conservation Biology*, 26(6):
494 1016–1025, 2012.
- 495 M. Chakraborty, A. Igarashi, W. Suksompong, and Y. Zick. Weighted envy-freeness in indivisible
496 item allocation. *TEAC*, 9(3):1–39, 2021.
- 497 Satya R. Chakravarty. *Ethical Social Index Numbers*. Springer Verlag, 1990.
- 498 Jingdi Chen, Yimeng Wang, and Tian Lan. Bringing fairness to actor-critic reinforcement learning
499 for network utility optimization. In *IEEE Conference on Computer Communications*, pp. 1–10,
500 2021.
- 501 Violet Xinying Chen and JN Hooker. A guide to formulating equity and fairness in an optimization
502 model. *Preprint*, pp. 162–174, 2021.
- 503 Yann Chevaleyre, Paul E Dunne, Michel Lemaître, Nicolas Maudet, Julian Padget, Steve Phelps,
504 and Juan A Rodríguez-aguilar. Issues in Multiagent Resource Allocation. *Computer*, 30:3–31,
505 2006.
- 506 Flavio Chierichetti, Ravi Kumar, Silvio Lattanzi, and Sergei Vassilvitskii. Fair clustering through
507 fairlets. *NeurIPS*, 30, 2017.
- 508 Alexandra Cimpeana, Catholijn Jonkerb, Pieter Libina, and Ann Nowéa. A multi-objective frame-
509 work for fair reinforcement learning. In *Multi-Objective Decision Making Workshop 2023*, 2023.
- 510 Cyrus Cousins. An axiomatic theory of provably-fair welfare-centric machine learning. In *NeurIPS*,
511 2021.
- 512 Virginie Do and Nicolas Usunier. Optimizing generalized gini indices for fairness in rankings. *arXiv*
513 *preprint arXiv:2204.06521*, 2022.
- 514 Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness
515 through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Con-*
516 *ference*, pp. 214–226, January 2012.
- 517 Zimeng Fan, Nianli Peng, Muhang Tian, and Brandon Fain. Welfare and fairness in multi-objective
518 reinforcement learning. *arXiv preprint arXiv:2212.01382*, 2022.
- 519 Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-
520 critic methods. In *ICML*, pp. 1582–1591, 2018.
- 521 Hoda Heidari, Claudio Ferrari, Krishna P. Gummadi, and Andreas Krause. Fairness behind a veil of
522 ignorance: A welfare analysis for automated decision making. In *NeurIPS*, 2018.

- 523 Shahin Jabbari, Matthew Joseph, Michael Kearns, Jamie Morgenstern, and Aaron Roth. Fairness in
524 reinforcement learning. In *ICML*, pp. 1617–1626, 2017.
- 525 N. Jensen. An introduction to bernoullian utility theory, I: utility functions. *Swedish Journal of*
526 *Economics*, 69:163–183, 1967.
- 527 Jiechuan Jiang and Zongqing Lu. Learning Fairness in Multi-Agent Systems. In *NeurIPS*, 2019.
- 528 Michael P Kim, Amirata Ghorbani, and James Zou. Multiaccuracy: Black-box post-processing for
529 fairness in classification. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and*
530 *Society*, pp. 247–254, 2019.
- 531 David Kurokawa, Ariel D. Procaccia, and Nisarg Shah. Leximin Allocations in the Real World.
532 In *Proceedings of the Sixteenth ACM Conference on Economics and Computation*, pp. 345–362,
533 June 2015. DOI: 10.1145/2764468.2764490.
- 534 Jurek Leonhardt, Avishek Anand, and Megha Khosla. User fairness in recommender systems. In
535 *Companion Proceedings of the The Web Conference 2018*, pp. 101–102, 2018.
- 536 Debmalaya Mandal and Jiarui Gan. Socially fair reinforcement learning. *arXiv preprint*
537 *arXiv:2208.12584*, 2022.
- 538 Dimitris Michailidis, Willem Röpke, Diederik M Roijers, Sennay Ghebreab, and Fernando P Santos.
539 Scalable multi-objective reinforcement learning with fairness guarantees using lorenz dominance.
540 *arXiv preprint arXiv:2411.18195*, 2024.
- 541 Kristof Van Moffaert and Ann Nowé. Multi-objective reinforcement learning using sets of pareto
542 dominating policies. *JMLR*, 15:3663–3692, 2014.
- 543 H. Moulin. *Fair Division and Collective Welfare*. MIT Press, 2004.
- 544 Yusuke Mukai, Yasuaki Kuroe, and Hitoshi Iima. Multi-objective reinforcement learning method
545 for acquiring all pareto optimal policies simultaneously. In *IEEE International Conference on*
546 *Systems, Man, and Cybernetics*, 2012.
- 547 Razieh Nabi, Daniel Malinsky, and Ilya Shpitser. Learning optimal fair policies. In *ICML*, 2019.
- 548 Samer B Nashed, Justin Svegliato, and Su Lin Blodgett. Fairness and sequential decision making:
549 Limits, lessons, and opportunities. *arXiv preprint arXiv:2301.05753*, 2023.
- 550 Patrice Perny, Paul Weng, Judy Goldsmith, and Josiah Hanna. Approximation of Lorenz-optimal so-
551 lutions in multiobjective Markov decision processes. In *International Conference on Uncertainty*
552 *in Artificial Intelligence (UAI)*, 2013.
- 553 M.L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. Wiley,
554 1994.
- 555 Junqi Qian, Umer Siddique, Guanbao Yu, and Paul Weng. From fair solutions to compromise
556 solutions in multi-objective deep reinforcement learning. *Neural Computing and Applications*,
557 pp. 1–31, 2025.
- 558 John Rawls. *The Theory of Justice*. Havard university press, 1971.
- 559 Mathieu Reymond, Eugenio Bargiacchi, and Ann Nowé. Pareto conditioned networks. *arXiv*
560 *preprint arXiv:2204.05036*, 2022.
- 561 Saeed Sharifi-Malvajerdi, Michael Kearns, and Aaron Roth. Average Individual Fairness: Algo-
562 rithms, Generalization and Experiments. In *NeurIPS*. 2019.

- 563 Umer Siddique, Paul Weng, and Matthieu Zimmer. Learning fair policies in multi-objective (deep)
564 reinforcement learning with average and discounted rewards. In *International Conference on*
565 *Machine Learning*, 2020.
- 566 Umer Siddique, Abhinav Sinha, and Yongcan Cao. Fairness in preference-based reinforcement
567 learning. In *ICML 2023 Workshop The Many Facets of Preference-Based Learning*, 2023.
- 568 Umer Siddique, Peilang Li, and Yongcan Cao. Fairness in traffic control: Decentralized multi-agent
569 reinforcement learning with generalized gini welfare functions. In *Multi-Agent reinforcement*
570 *Learning for Transportation Autonomy*, 2024a.
- 571 Umer Siddique, Peilang Li, and Yongcan Cao. Towards fair and equitable policy learning in co-
572 operative multi-agent reinforcement learning. In *Coordination and Cooperation for Multi-Agent*
573 *Reinforcement Learning Methods Workshop*, 2024b.
- 574 Ashudeep Singh and Thorsten Joachims. Policy Learning for Fairness in Ranking. In *NeurIPS*.
575 2019.
- 576 Till Speicher, Hoda Heidari, Nina Grgic-Hlaca, Krishna P. Gummadi, Adish Singla, Adrian Weller,
577 and Muhammad Bilal Zafar. A unified approach to quantifying algorithmic unfairness: Measuring
578 individual & group unfairness via inequality indices. In *ACM SIGKDD International Conference*
579 *on Knowledge Discovery and Data Mining*, 2018.
- 580 Kristof Van Moffaert and Ann Nowé. Multi-objective reinforcement learning using sets of pareto
581 dominating policies. *The Journal of Machine Learning Research*, 15(1):3483–3512, 2014.
- 582 Kristof Van Moffaert, Madalina M Drugan, and Ann Nowé. Scalarized multi-objective reinforce-
583 ment learning: Novel design techniques. In *IEEE Symposium on Adaptive Dynamic Programming*
584 *and Reinforcement Learning (ADPRL)*, pp. 191–199, 2013.
- 585 Paul Weng. Fairness in reinforcement learning. In *AI for Social Good Workshop at International*
586 *Joint Conference on Artificial Intelligence*, 2019.
- 587 J.A. Weymark. Generalized Gini inequality indices. *Mathematical Social Sciences*, 1:409–430,
588 1981.
- 589 R.R. Yager. On ordered weighted averaging aggregation operators in multi-criteria decision making.
590 *IEEE Trans. on Syst., Man and Cyb.*, 18:183–190, 1988.
- 591 Runzhe Yang, Xingyuan Sun, and Karthik Narasimhan. A generalized algorithm for multi-objective
592 reinforcement learning and policy adaptation. *NeurIPS*, 32, 2019.
- 593 Guanbao Yu, Umer Siddique, and Paul Weng. Fair deep reinforcement learning with generalized
594 gini welfare functions. In *Adaptive and Learning Agents (ALA) Workshop*, 2023a.
- 595 Guanbao Yu, Umer Siddique, and Paul Weng. Fair deep reinforcement learning with preferential
596 treatment. In *ECAI*, 2023b.
- 597 Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, Krishna P. Gummadi, and Adrian
598 Weller. From Parity to Preference-based Notions of Fairness in Classification. In *NIPS*, 2017.
- 599 Xueru Zhang and Mingyan Liu. Fairness in learning-based sequential decision algorithms: A survey.
600 In *Handbook of Reinforcement Learning and Control*, pp. 525–555. Springer, 2021.
- 601 Matthieu Zimmer, Claire Glanois, Umer Siddique, and Paul Weng. Learning fair policies in decen-
602 tralized cooperative multi-agent reinforcement learning. In *ICML*, 2021.