
Envy-free Policy Teaching to Multiple Agents

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study envy-free policy teaching. A number of agents independently explore a
2 common Markov decision process (MDP), but each with their own reward function
3 and discounting rate. A teacher wants to teach a target policy to the diverse group
4 of agents, by way of modifying the agents' reward functions, providing additional
5 bonus to certain behaviors or penalizing others. These reward modifications are
6 personalized for each agent. An important question in this setting concerns how
7 a teaching program can be designed so that the agents think that they are treated
8 fairly. We adopt the fairness notion of envy-freeness (EF) to formalize this question
9 and define three different EF notions, each imposing stronger requirements than
10 the previous one. Using these notions, we then investigate several fundamental
11 questions, including the existence of EF solutions in the policy teaching setting,
12 the computation of cost-minimizing solutions, and the price of fairness (PoF), i.e.,
13 the increase in cost due to consideration of fairness. We show that an EF solution
14 may not exist when penalties are not allowed, but exists otherwise. Depending on
15 the cost measures, computing a cost-minimizing EF solution can be formulated as
16 convex or linear programming and hence solved efficiently. Asymptotically, the
17 PoF increases but at most linearly with the geometric sum of the discount factor in
18 general, the size of the MDP, and the number of agents involved. Thus, fairness
19 can be incorporated in multi-agent teaching without significant computational or
20 price-of-fairness burdens.

21 1 Introduction

22 Incentive design is an important approach to influencing rational agents' behavior. In reinforcement
23 learning (RL), the incentive of an agent is expressed through their reward function, which plays a key
24 role in determining the policy the agent learns [1]. One can thus teach a desired policy to an agent by
25 modifying the agent's reward function, in a way that makes the target policy optimal with respect
26 to the new rewards. This can be very useful in many scenarios. In safe RL, for example, penalties
27 can be imposed on dangerous actions to prevent an agent from executing them [2]. In many cases,
28 personalized teaching programs are useful when agents to be taught are heterogeneous: they might
29 have very different innate reward functions, or apply different discount factors when evaluating a
30 policy. As a result, the agents may find them rewarded/penalized differently for performing the same
31 action in the same situation (see Figure 1). Concerns of fairness may arise due to these differences,
32 leading us to the question of how to design personalized teaching programs so that agents being
33 taught think that they are treated fairly.

34 For example, in a classroom teaching setting, a teacher wants to teach a group of students to
35 accomplish a task that involves performing a sequence of actions. Personalize teaching schemes
36 might be used concerning the diverse background of the students, and it is usually important that the
37 students feel they are treated equally by the teacher. Similarly, in a principal-agent setting, a company
38 outsources a task to different contractors. Rewards or penalties are stipulated via personalized

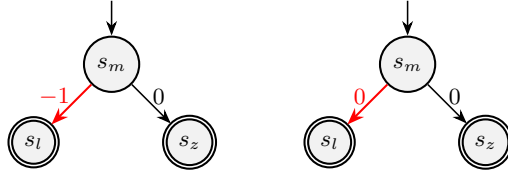


Figure 1: To teach the agent to choose the action leading to state s_l , an additional reward 1 is necessary for an agent whose innate reward function is illustrated on the left, whereas an agent with the reward function on the right already find this target policy optimally. When the personalized teaching programs are used, the agent on the right may think they are treated unfairly as they get no bonus for following the target policy while the agent on the left gets bonus 1.

39 contracts to ensure that the contractors abide by a desired policy when fulfilling the task. Meanwhile
 40 fairness might also be an important consideration as a beneficial factor for long-term partnerships.

41 Towards answering this question, the first step is to understand what it means to be fair in the setting
 42 of policy teaching. Indeed, in a world with growing awareness of equality and transparency, fairness
 43 has discussed and evaluated in a wide range of applications. Various fairness notions and concepts
 44 have been proposed and used across many domains in technology [3]. We borrow the well-studied
 45 fairness notion of *envy-freeness* (EF) from the literature on fair division, which naturally applies to
 46 the policy teaching setting. The notion is widely used for settling disputes over property divisions,
 47 or finding resource allocations that make participants with different valuations unenvious of each
 48 other (e.g., the fair cake-cutting problem [4]). By adapting EF to policy teaching, we aim to find a set
 49 of personalized teaching programs for a group of agents, such that no agent would prefer to switch
 50 the program they receive with another agent. At the same time, as a basic requirement of policy
 51 teaching, each program should also incentivize the corresponding agent who receives it to use the
 52 target policy. Besides this basic notion, we also consider two stronger variants of it—one allows an
 53 agent to further deviate from the target policy when evaluating the potential benefit had they received
 54 another agent’s teaching program; the other simply requires all teaching programs to be identical,
 55 which gives complete fairness in a sense.

56 **Main Results** We investigate several fundamental questions related to the above laid out objective.

- 57 • *Existence of a Solution.* The first question is about the existence of an EF solution under the
 58 three EF notions of interest. We show that an EF solution always exists and one can be obtained
 59 simply by penalizing undesired actions by a sufficiently large value. Nevertheless, the reverse does
 60 not hold true: one cannot hope to find an EF solution by rewarding actions desired by the target
 61 policy, no matter how large the rewards are set. We demonstrate instances that do not admit any EF
 62 solution even with the weakest EF notion when only bonuses are allowed, but also prove that this
 63 non-existence issue is resolved if the agents have an identical discount factor.
- 64 • *Cost Minimization.* Since reward modification can be very costly, we are also interested in finding
 65 out a solution with the least cost. We consider two cost measures in this paper, which consider
 66 the norm of the reward modifications made and the cumulative cost the teacher actually pays,
 67 respectively. Depending on the choice of the cost measure, we show that computing a cost-
 68 minimizing EF solution can be formulated as convex or linear programming and can hence be
 69 solved efficiently.
- 70 • *Price of Fairness.* Finally, we analyze the price of fairness (PoF), a quantity that measures the
 71 (multiplicative) increase of cost due to consideration of fairness—in a similar spirit of the PoA
 72 (price of anarchy) in game theory. We show that, asymptotically, the PoF increases linearly with the
 73 geometric sum of the discount factor in general, while it may also grow, at most linearly, with the
 74 size of the MDP and the number of agents involved depending on the specific EF notion considered.

75 In summary, our results indicate that the consideration of fairness in addition to the original goal of
 76 policy teaching may result in non-existence of workable solutions, but the existence is guaranteed
 77 in a fairly wide range of important settings. It does not appear to complicate the problem in terms
 78 of computational complexity. The additional cost it incurs grows moderately with the size of the

79 problem. Thus, fairness can be incorporated in multi-agent teaching without significant computational
80 or price-of-fairness burdens.

81 **Related Work** Our work lies at the intersection of policy teaching and envy-free resource allocation.

82 Without the fairness constraints, our model can be seen as a policy teaching problem for each
83 individual agent in the model. A number of studies have looked at this problem [5, 6]. The problem
84 can be computationally harder though when the target is to hit one in a set of policies rather than
85 a single target [7]. When the teacher is targeting a malicious policy, policy teaching can also be
86 interpreted as reward poisoning [8, 9, 10, 11, 12]. From a technical point of view, these two problems
87 are almost identical and can be solved by using the same techniques. However, conceptually, it is
88 less likely that one would take fairness into consideration when designing a poisoning attack. More
89 broadly, policy teaching can be seen as a sub-field of reward design, a broader area that studies how
90 to influence agents’ behaviors thought tweaking the reward function. The objectives of these studies
91 are not limited to inducing a target policy. A notable example is *reward shaping* [13, 14, 15], which
92 aims to accelerate an agent’s learning process through reward design. Indeed, while our focus is on
93 policy teaching, the same question of how to design rewards fairly can be asked with other objectives
94 as well. These can be potential directions for future work.

95 The study of fair division dates back to the early work of Foley [16], and the formal concept of
96 envy-freeness appeared even earlier [17]. Research on fair division has since evolved into a large
97 body of work, with focuses on allocation of divisible or indivisible items [18, 19, 20]. Our work is in
98 particular related to fair allocation of indivisible goods with subsidies [21], as external benefits are
99 provided to change the agents’ original incentives in both settings. Our work considers the additional
100 goal of teaching the target policy; indeed, without this goal, an envy-free solution can be achieved
101 trivially by no making any modification to the original reward functions.

102 We note that there are also other studies on machine teaching settings with multiple agents, or
103 multiple teachers [22, 23, 24], though with very different models from ours. From a mechanism
104 design perspective, our model can also be viewed as one version of the contract design problem [25],
105 where a principal offers an agent a contract for performing a target policy, but might be uncertain
106 about the agent’s type (i.e., the original reward function). Our EF solutions correspond exactly to
107 truthful mechanisms that elicit the agent’s true type.

108 2 Preliminaries

109 There is a set of n agents $1, \dots, n$. Let $[n] = \{1, \dots, n\}$. Each agent $i \in [n]$ faces an MDP
110 $\mathcal{M}_i = \langle S, A, R_i, P, \mathbf{z}, \gamma_i \rangle$ with shared state space S , action space A , transition function $P : S \times A \times$
111 $S \rightarrow [0, 1]$ describing the state transition dynamics, and initial state distribution \mathbf{z} . Moreover, there is
112 a personal reward function $R_i : S \times A \rightarrow \mathbb{R}$ and discount factor γ_i . Whenever agent i takes an action
113 a in state s , a reward $R_i(s, a)$ is generated for this agent; meanwhile the state transitions to a next
114 one $s' \in S$ with probability $P(s, a, s')$. We consider the setting where each agent is concerned with
115 the (expected) cumulative reward, i.e., the discounted sum of rewards with respect to γ_i , obtained
116 over an infinite horizon. Without loss of generality these are actions generated by following some
117 deterministic policy $\pi : S \rightarrow A$, and the cumulative reward of this policy for agent i is

$$\rho_i^\pi := \mathbb{E} [\sum_{t=0}^{\infty} (\gamma_i)^t \cdot R_i(s_t, a_t) | s_0 \sim \mathbf{z}, \pi],$$

118 where the expectation is taken over the trajectory $(s_t, a_t)_{t=0}^{\infty}$ resulting from an initial state s_0 sampled
119 from \mathbf{z} and the agent executing π subsequently. The agent aims to find an optimal policy, which
120 maximizes ρ_i^π , and this can usually be handled by standard planning and reinforcement learning
121 algorithms.

122 Note that in the setting we consider throughout this paper, the agents operate independently in
123 separate environments.

124 2.1 Single-agent Policy Teaching

125 Consider the situation where we want agent i to execute a target policy π^* , but the agent finds a
126 different policy π' optimal for \mathcal{M}_i . We hope to influence the agent’s decision making and incentivize
127 them to use π^* . A typical way is by modifying the agent’s reward function by providing additional

128 rewards. We follow the literature and consider only deterministic target policies—indeed, in general,
 129 one cannot hope to incentivize an agent to use precisely a non-deterministic policies only through
 130 tweaking the reward function.

131 Specifically, the teacher chooses a *reward adjustment function* $\delta_i : S \times A \rightarrow \mathbb{R}$, or *adjustment* for
 132 short, whereby an additional reward $\delta_i(s, a)$ is provided whenever the agent takes an action $a \in A$
 133 in a state $s \in S$. Effectively, the adjustment changes the agent’s reward function to $\tilde{R}_i(s, a) =$
 134 $R_i(s, a) + \delta_i(s, a)$. The agent then optimizes their policy with respect to \tilde{R}_i , and will be incentivized
 135 to use π^* if it offers the maximum *cumulative reward* with respect to \tilde{R}_i . The cumulative reward of
 136 agent i obtained by executing a policy π is defined as the sum of discounted rewards as follows;
 137 we view it as a function of δ_i :

$$\rho_i^\pi(\delta_i) := \mathbb{E} \left[\sum_{t=0}^{\infty} (\gamma_i)^t \cdot \tilde{R}_i(s_t, a_t) \mid s_0 \sim \mathbf{z}, \pi \right].$$

138 We can define the V-function and Q-function of π conditioned on adjustment δ_i as follows. For all
 139 $s \in S$ and $a \in A$:

$$V_i^\pi(s|\delta_i) = Q_i^\pi(s, \pi(s)|\delta_i),$$

and $Q_i^\pi(s, a|\delta_i) = \tilde{R}_i(s, a) + \gamma_i \cdot \mathbb{E}_{s' \sim P(s, a, \cdot)} V_i^\pi(s'|\delta_i).$

140 The V-function captures the cumulative reward the agent obtains by starting from s and following
 141 π , whereas the Q-function captures the cumulative reward by starting from s and taking action a
 142 at the first step. We have $\rho_i^\pi(\delta_i) = V_i^\pi(\mathbf{z}|\delta_i) := \mathbb{E}_{s_0 \sim \mathbf{z}} V_i^\pi(s_0|\delta_i)$. Using these two concepts, the
 143 Bellman equation further gives the optimal policy in the MDP: a policy π is optimal if and only if
 144 $Q_i^\pi(s, \pi(s)|\delta_i) \geq Q_i^\pi(s, a|\delta_i)$ for all $s \in S$ and $a \in A$.

145 **Incentive Constraint** Hence, we require that π^* satisfies the condition described by the Bellman
 146 equation. Since the agent may find multiple policies optimal, a robustness guarantee $\epsilon > 0$ is
 147 considered to strictly incentivize the agent to use π^* . The incentive constraint requires that

$$Q_i^{\pi^*}(s, \pi^*(s)|\delta_i) \geq Q_i^{\pi^*}(s, a|\delta_i) + \epsilon \quad \text{for all } a \neq \pi^*(s),$$

148 so that π^* remains optimal even if there is a small noise in the agent’s perception of the values.

149 **Cost Measures** The goal of the teacher is to find the most cost-efficient way of teaching. In
 150 particular, we consider the following cost measures in this paper:

$$\text{cost}(\delta_i) = \|\delta_i\| = \left(\sum_{s \in S, a \in A} (\delta_i(s, a))^2 \right)^{1/2}, \quad (1)$$

151 which measures the norm of the adjustment; or

$$\text{cost}(\delta_i) = \mathbb{E} \left[\sum_{t=0}^{\infty} \bar{\gamma}^t \cdot \delta_i(s_t, a_t) \mid s_0 \sim \mathbf{z}, \pi \right], \quad (2)$$

152 which measures the cumulative payment, where $\bar{\gamma}$ is a discount factor from the teacher’s perspective.
 153 The cumulative payment measure applies in particular when the teacher uses non-negative rewards
 154 (i.e., subsidies) as additional incentives for the agent.

155 3 Teaching Multiple Agents and Fairness

156 In the multi-agent setting, the teacher needs to provide adjustments to all the agents in $[n]$. We call
 157 a collection of adjustments $(\delta_i)_{i \in [n]}$ an *adjustment scheme*. A basic approach for this setting is to
 158 deal with each agent separately, by solving a single-agent teaching problem for each of them. The
 159 solution obtained via this approach provides personalized adjustments to the agents and it minimizes
 160 the teacher’s overall cost. Nevertheless, it might not be fair as we showed in the example of Figure 1
 161 To be more specific, we define three fairness notions, each being stronger than the previous one. We
 162 start with the following weak EF notion.

163 **Definition 3.1 (Weak envy-freeness (WEF)).** An adjustment scheme $(\delta_i)_{i \in [n]}$ is *envy-free* if it holds
 164 for all $i \in [n]$ that

$$\rho_i^{\pi^*}(\delta_i) \geq \rho_i^{\pi^*}(\delta_j) \quad \text{for all } j \in [n]. \quad (3)$$

165 In other words, no agent i would prefer the adjustment for another agent j to their own.

166 The above notion only compares the agents' benefits under π^* . Indeed, we require that δ_i incentivizes
 167 agent i to use π^* , so the left side of (3) is exactly also the highest possible extra benefit i can obtain
 168 given adjustment δ_i . But this is not true for the adjustment on the right side: π^* need not be optimal for
 169 agent i with respect to $R_i + \delta_j$; a higher cumulative reward might be attainable if the agent switches
 170 to another policy. In some scenarios, this higher potential reward may be a legitimate consideration
 171 when fairness is evaluated. The following stronger notion takes this aspect into consideration.

172 **Definition 3.2 (Envy-freeness (EF)).** An adjustment scheme $(\delta_i)_{i \in [n]}$ is *strongly envy-free* if it holds
 173 for all $i \in [n]$ that:

$$\max_{\pi} \rho_i^{\pi}(\delta_i) \geq \max_{\pi} \rho_i^{\pi}(\delta_j) \quad \text{for all } j \in [n]. \quad (4)$$

174 An even stronger fairness notion defined below simply requires the same adjustment to be applied to
 175 all the agents.

176 **Definition 3.3 (Strong envy-freeness (SEF)).** An adjustment scheme $(\delta_i)_{i \in [n]}$ is *completely fair* if
 177 $\delta_1 = \dots = \delta_n$.

178 Let \mathcal{D}_{WEF} , \mathcal{D}_{EF} , and \mathcal{D}_{SEF} denote the sets of adjustment schemes complying with the above fairness
 179 notions respectively. It is not hard to see that: $\mathcal{D}_{\text{WEF}} \supseteq \mathcal{D}_{\text{EF}} \supseteq \mathcal{D}_{\text{SEF}}$.

180 Besides the fairness guarantees, we also want the adjustment scheme to—as our initial objective—
 181 incentivize the agents to use the target policy, i.e., to satisfy the incentive constraints we defined above;
 182 we call such adjustment schemes *feasible* schemes (Definition 3.4). Indeed, without the feasibility
 183 guarantee, all the fairness notions above can be trivially satisfied by providing zero additional reward
 184 to every agent. In addition, sometimes only non-negative additional rewards are allowed, e.g., when
 185 we can provide the agents with subsidies but cannot penalize them. Hence, we define *non-negative*
 186 adjustment schemes (Definition 3.5).

187 **Definition 3.4 (Feasibility).** An adjustment scheme $(\delta_i)_{i \in [n]}$ is *feasible* (with respect to a given
 188 robustness guarantee $\epsilon > 0$) if for all $i \in [n]$, it holds that

$$\tilde{Q}_i^{\pi^*}(s, \pi^*(s)) \geq \tilde{Q}_i^{\pi^*}(s, a) + \epsilon \quad \text{for all } a \neq \pi^*(s). \quad (5)$$

189 **Definition 3.5 (Non-negativity).** An adjustment scheme $(\delta_i)_{i \in [n]}$ is *non-negative* if $\delta_i(s, a) \geq 0$ for
 190 all $i \in [n]$, $s \in S$, and $a \in A$.

191 Similarly to the single-agent teaching problem, we are interested in obtaining a cost-minimizing
 192 solution. We consider the sum of the costs and define the cost of an adjustment scheme δ as

$$\text{cost}(\delta) = \sum_{i \in [n]} \text{cost}(\delta_i).$$

193 Next, before we delve into the computation of a cost-minimizing solution, we first investigate the
 194 existence of a solution with respect to the above defined fairness notions.

195 4 Existence of a Fair Feasible Solution

196 Throughout this section, we assume that the original rewards are bounded in the interval $[-h, h]$, i.e.,
 197 $R_i(s, a) \in [-h, h]$ for all s, a , and i . Our first result shows that a fair and feasible solution always
 198 exists under all of the above fairness notions, in particular under the strongest notion SEF.

199 **Theorem 4.1.** *For any robustness guarantee $\epsilon > 0$, an SEF and feasible adjustment scheme always*
 200 *exists*¹

201 *Proof sketch.* The idea is to penalize actions off the target policy by a sufficiently large value. We
 202 construct an adjustment scheme $(\delta_i)_{i \in [n]}$ where

$$\delta_i(s, a) = \begin{cases} 0, & \text{if } a = \pi^*(s) \\ -\max_{i' \in [n]} \frac{2h}{1-\gamma_{i'}} - \epsilon, & \text{otherwise} \end{cases}$$

¹Full proofs and omitted proofs can all be found in the appendix.

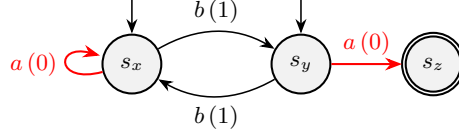


Figure 2: There are two agents, whose discount factors are $\gamma_1 = 0.9$ and $\gamma_2 = 0.5$, respectively. $S = \{s_x, s_y, s_z\}$, $A = \{a, b\}$, and all transitions are deterministic. Originally, the agents' reward functions are the same, with $R_1(s, a) = R_2(s, a) = 0$ and $R_1(s, b) = R_2(s, b) = 1$ for all $s \in S$, as annotated on the edges. The states s_x and s_y are chosen as the initial state with equal probability. The target policy π^* , highlighted as the red edges, is such that $\pi^*(s) = a$ for all $s \in S$, i.e., it always chooses action a .

203 for all $s \in S$ and $i \in [n]$. The scheme is obviously SEF as δ_i is the same for all the agents. It can also
 204 be verified that it is feasible. Intuitively, the penalty is so large such that once the agent is penalized,
 205 the subsequent cumulative rewards cannot compensate this penalty even if the highest rewards are
 206 attained at every step. \square

207 Nevertheless, if we only allow non-negative schemes, the existence of a feasible solution cannot be
 208 taken for granted. As we prove in Theorem 4.2, the example illustrated in Figure 2, albeit involving
 209 only two agents, does not admit any EF feasible solution (and hence neither an SEF one).

210 **Theorem 4.2.** *For any robust guarantee $\epsilon \geq 0$, a feasible adjustment scheme that is WEF and*
 211 *non-negative may not exist, even when there are only two agents and the agents' reward functions are*
 212 *the same (but discount factors are different).*

213 *Proof.* We show that there exists no feasible adjustment scheme that is WEF and non-negative in
 214 the example illustrated in Figure 2. Suppose for the sake of contradiction that there exists a scheme
 215 (δ_1, δ_2) which is EF, non-negative, and feasible.

216 Without loss of generality, we can assume that $\delta_1(s, b) = \delta_2(s, b) = 0$ for all $s \in S$. Indeed, it is not
 217 hard to see that if there exists a WEF and feasible scheme with some or all of these values being
 218 strictly positive, it will remain WEF and feasible if these values are reset to 0. Hence, it remains to
 219 pin down the values for action a in the adjustment scheme. For ease of description, let $x_i = \delta_i(s_x, a)$
 220 and $y_i = \delta_i(s_y, a)$ for $i \in \{1, 2\}$.

221 We first argue that the following two inequalities hold:

$$y_1 \geq y_2, \quad \text{and} \quad x_2 \geq x_1. \quad (6)$$

222 To see this, consider the WEF constraints. The adjustment scheme considered is WEF, so it satisfies
 223 (3), which implies that $\rho_i^{\pi^*}(\delta_i) \geq \rho_i^{\pi^*}(\delta_{-i})$, where $-i$ means the other index in $\{1, 2\}$. Hence,

$$0.5 \cdot V_i^{\pi^*}(s_x | \delta_i) + 0.5 \cdot V_i^{\pi^*}(s_y | \delta_i) \geq 0.5 \cdot V_i^{\pi^*}(s_x | \delta_{-i}) + 0.5 \cdot V_i^{\pi^*}(s_y | \delta_{-i}). \quad (7)$$

224 It is easy to derive the values of s_x and s_y as neither of them depends on the values of any other
 225 states. When agent i uses π^* , we have

$$V_i^{\pi^*}(s_x | \delta_j) = Q_i^{\pi^*}(s_x, a | \delta_j) = \frac{1}{1-\gamma_i} \cdot y_j, \quad (8)$$

$$\text{and} \quad V_i^{\pi^*}(s_y | \delta_j) = Q_i^{\pi^*}(s_y, a | \delta_j) = x_j. \quad (9)$$

226 Plugging these two equations back into (7) gives

$$0.5 \cdot \frac{1}{1-\gamma_i} \cdot y_i + 0.5 \cdot x_i \geq 0.5 \cdot \frac{1}{1-\gamma_{-i}} \cdot y_{-i} + 0.5 \cdot x_{-i},$$

227 where 0.5 is the probability of the initial states. Replacing γ_i with the corresponding values gives

$$10 \cdot y_1 + x_1 \geq 10 \cdot y_2 + x_2 \quad (10)$$

$$2 \cdot y_2 + x_2 \geq 2 \cdot y_1 + x_1 \quad (11)$$

228 It follows that (10) + (11) gives $y_1 \geq y_2$, and (10) + 5 · (11) gives $x_2 \geq x_1$.

229 Next, we turn to the feasibility constraints. The assumption that δ_i is feasible means that

$$Q_i^{\pi^*}(s_x, a|\delta_i) \geq Q_i^{\pi^*}(s_x, b|\delta_i) + \epsilon = 1 + \gamma_i \cdot \tilde{V}_i^{\pi^*}(s_y) + \epsilon$$

and $Q_i^{\pi^*}(s_y, a|\delta_i) \geq Q_i^{\pi^*}(s_y, b|\delta_i) + \epsilon = 1 + \gamma_i \cdot \tilde{V}_i^{\pi^*}(s_x) + \epsilon$

230 Plugging the value function and Q-function ((8) and (9)) into the above two equations gives:

$$1 + \frac{\gamma_i}{1-\gamma_i} \cdot y_i \leq x_i \leq \frac{1}{\gamma_i(1-\gamma_i)} \cdot y_i - \frac{1}{\gamma_i}. \quad (12)$$

231 Hence, applying (12) and (6), we have

$$9 \cdot y_1 + 1 \leq x_1 \leq x_2 \leq 4 \leq 4 \cdot y_1 - 2,$$

232 This means that $y_1 < 0$ and contradicts the assumption that δ is a non-negative scheme. \square

233 It turns out that the agents' discount factors also play a crucial role: an identical discount factor is
 234 sufficient for ensuring the existence of a feasible SEF solution. We present this result below but leave
 235 the proof to the appendix.

236 **Theorem 4.3.** *When the agents have an identical discount factor, a feasible adjustment scheme that*
 237 *is also SEF and non-negative always exists, for any robustness guarantee $\epsilon > 0$.*

238 *Proof sketch.* Suppose that $\gamma_1 = \dots = \gamma_n = \gamma$. Let $H = \frac{2\cdot\gamma}{1-\gamma} \cdot h + 2h + \epsilon$. We construct the
 239 following scheme $\delta = (\delta_i)_{i \in [n]}$:

$$\delta_i(s, a) = \begin{cases} H \cdot \sum_{s' \in S \setminus S^T} P(s, a, s') + \frac{1}{1-\gamma} \cdot H \cdot \sum_{s' \in S^T} P(s, a, s'), & \text{if } a = \pi^*(s) \\ 0, & \text{otherwise} \end{cases} \quad (13)$$

240 for all $s \in S$ and $i \in [n]$, where S^T denotes the set of terminal states in S .

241 The scheme is obviously non-negative and SEF. Moreover, we can also show that it is feasible.
 242 Intuitively, δ_i results the agent getting a reward that is sufficiently large (and is roughly the same) at
 243 every step if they follow π^* . The rewards are adjusted by a factor of $1/(1-\gamma)$ at terminal states so
 244 that it is as if the process continues forever. Hence, under δ_i , the process is roughly the same as an
 245 infinite process where the agent gets a constant positive reward at every step if they take an action
 246 stipulated by π^* , and reward 0 otherwise—the optimal choice for the agent in such a process is to
 247 always follow π^* . \square

248 5 Computing an Optimal Fair Solution

249 In this section, we investigate the computational problems of fair policy teaching. Our main result is
 250 that for each of the fair notions we defined the set of fair solutions lie in a convex polytope defined by
 251 polynomially many linear constraints. To find out a cost minimizing solution from this convex set
 252 can then be formulated as a convex optimization problem, given that the two types of cost functions
 253 we consider (i.e., (1) and (2)) are convex functions of the adjustment scheme. We show how to
 254 write down the various types of constraints considered as linear constraints next. We start from the
 255 feasibility constraints.

256 **Feasibility Constraints** A feasible scheme δ is characterized by the following linear constraints,
 257 where we add an auxiliary variable $V_i(s)$ for each $s \in S$, and $Q_i(s, a)$ for pair $(s, a) \in S \times A$ —
 258 these auxiliary variables correspond to the V- and Q-functions of the target policy when δ is applied.

$$V_i(s) = Q_i(s, \pi^*(s)) \quad \forall i, s \quad (14a)$$

$$Q_i(s, a) = R_i(s, a) + \delta_i(s, a) + \gamma_i \sum_{s' \in S} P(s, a, s') \cdot V_i(s') \quad \forall i, s, a \quad (14b)$$

$$Q_i(s, \pi^*(s)) \geq Q_i(s, a) + \epsilon \quad \forall i, s, a \neq \pi^*(s) \quad (14c)$$

259 Indeed, the first two lines above follow from the Bellman equation and captures the values $V_i^{\pi^*}(s|\delta_i)$
 260 and $Q_i^{\pi^*}(s, a|\delta_i)$; the last line is the incentive constraints and enforces the feasibility of δ .

261 We then consider each of the fairness notions.

| | PoWEF | PoEF | PoSEF |
|--------------------------------|----------------------------------|--|--|
| Norm (1) | $\Theta(\lambda \cdot \sqrt{m})$ | $\Theta(\lambda \cdot n \cdot \sqrt{m})$ | $\Theta(\lambda \cdot n \cdot \sqrt{m})$ |
| Cumulative (2) | $\Theta(\lambda \cdot n)$ | $\Theta(\lambda \cdot n)$ | $\Theta(\lambda \cdot n)$ |

Table 1: Summary of the price of fairness.

262 **SEF Constraints** To enforce SEF simply amounts to the following constraints for each pair of
263 agents $i, j \in [n]$, which enforces the schemes to be identical:

$$\delta_i(s, a) = \delta_j(s, a) \quad \forall s, a. \quad (15)$$

264 **WEF Constraints** To enforce WEF, we add additional variables $V_{i,j}$ and $Q_{i,j}$ to capture the
265 values $V_i^{\pi^*}(s|\delta_j)$ and $Q_i^{\pi^*}(s, a|\delta_j)$, i.e., the values agent i would have got had they been offered the
266 adjustment for agent j . We then add the following constraints—the same format as the Bellman
267 equation—to force these additional variables to acquire the desired values.

$$V_{i,j}(s) = Q_{i,j}(s, \pi^*(s)) \quad \forall i, j, s \quad (16a)$$

$$Q_{i,j}(s, a) = R_i(s, a) + \delta_j(s, a) + \gamma_i \sum_{s' \in S} P(s, a, s') \cdot V_{i,j}(s') \quad \forall i, j, s, a \quad (16b)$$

268 Thus, WEF simply amounts to the following constraints for each pair of agents $i, j \in [n]$ (recall that
269 \mathbf{z} is the distribution of the initial state):

$$\sum_{s \in S} z_s \cdot V_i(s) \geq \sum_{s \in S} z_s \cdot V_{i,j}(s) \quad \forall s, a \quad (17)$$

270 **EF Constraints** Similarly to the approach for handling the WEF constraints, we need additional
271 variables to capture the values of each agent i had then been offered adjustment δ_j . Indeed, we also
272 use Constraints in [\(16\)](#) and [\(17\)](#) but replace [\(16a\)](#) with the following one:

$$V_{i,j}(s) \geq Q_{i,j}(s, a) \quad \forall i, j, s, a \quad (18a)$$

273 which associates $V_{i,j}(s)$ to the maximum $Q_{i,j}(s, a)$, instead of $Q_{i,j}(s, \pi^*(s))$.

274 It should be noted that under these constraints, $V_{i,j}(s)$ in a solution will not necessarily be equal to
275 the V-values of the optimal policy—it will only be an upper bound of them. However, this will not
276 cause any issue to our approach as our goal is not to compute the V-values. Indeed, if the actual
277 V-values does not satisfy [\(17\)](#), then these constraints cannot be satisfied by any upper bounds of the
278 actual V-values, either.

279 **Non-negativity Constraints** Finally, to enforce non-negativity, we simply need the additional
280 constraints: $\delta_i(s, a) \geq 0$ for all i, s , and a .

281 6 Price of Fairness

282 We consider the price of fairness (PoF) in this section. The PoF measures the increase of teaching
283 cost due to consideration of fairness. It is in a similar spirit of the well-known concept of the price
284 of anarchy in economics and game theory. More specifically, we define the PoWEF, PoEF, and
285 PoSEF for our three fairness notions, which stand for the prices of WEF, EF, and SEF, respectively.
286 Formally, let $\mathcal{I}_{n,m,\lambda}$ be the set of instances with n agents, m state-action pairs (i.e., $m = |S| \cdot |A|$),
287 and $\frac{1}{1-\gamma_i} \leq \lambda$ for all $i \in [n]$. We define

$$\text{PoEF}(n, m, \lambda) := \max_{I \in \mathcal{I}_{n,m,\lambda}} \frac{\min_{\delta: \text{EF and feasible for } I} \text{cost}(\delta)}{\min_{\delta: \text{feasible for } I} \text{cost}(\delta)}.$$

288 PoWEF and PoSEF can be defined the same way with the corresponding notions. We consider
289 both types of cost functions we defined in Section [2](#). When cumulative payment is considered (i.e.,
290 [\(2\)](#)), we require adjustment schemes to be non-negative for this measure to be meaningful. Since a
291 feasible fair solution may not exist with this requirement, we only consider instances with an identical
292 discount factor for all the agents, in which case the existence of a feasible fair solution is guaranteed
293 (see Theorem [4.3](#)).

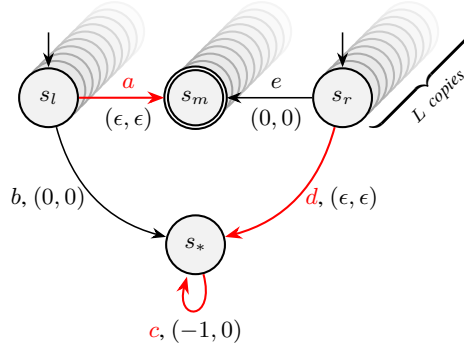


Figure 3: There are n agents with discount factors $\gamma_1 = \dots = \gamma_n = \gamma$. $A = \{a, b, c, d\}$ and all transitions are deterministic. The original rewards for agents 1 and 2 are annotated on the corresponding edges (numbers in parentheses), which are the same for both agents except the ones for c . Moreover, the reward functions are completely identical for agents $2, \dots, n$. There are L sets of copies of s_l , s_m , and s_r , each connected to copies in the same set and in addition s_* in the same way illustrated above. The initial state follows a uniform distribution over s_l , s_m , and all their copies. The target policy is highlighted in red: $\pi^*(s_l) = a$, $\pi^*(s_r) = d$, and $\pi^*(s_*) = c$.

294 We analyze the asymptotic growth of the PoF. The results are summarized in Table 1. The PoF
 295 increases linearly with λ in all the cases, and with the number of agents involved, with the only
 296 exception of PoWEF when the norm is used as the cost function. It also increases but sub-linearly
 297 with the size of the MDP with this cost function but not with the cumulative payment. Our results are
 298 all asymptotically tight.

299 Due to space limit, we leave the detailed proofs of the bounds to the appendix and only provide
 300 some intuition about the growth of PoF here. Consider the example illustrated in Figure 3. Without
 301 fairness consideration, all agents except agent 1 already find the target policy optimal, whereas agent
 302 1 prefers action e to d at state s_r . Overall, it suffices to give a bonus of 1 on (s_*, c) for agent 1, the
 303 cost of which is 1. Now we consider fairness constraints, and suppose that we still provide a bonus
 304 $\delta_1(s_*, c) = 1$. The consequence is that agents $2, \dots, n$ will be envious of this bonus to agent 1. To
 305 achieve SEF for example, the same bonus will have to be offered to these agents as well. Nevertheless,
 306 while this removes envies, it also incentivizes the agents to take action b instead of a , leading to
 307 violation of the feasibility constraint. Inevitably, to construct a feasible and fair in this example, we
 308 cannot hope to only modify the reward for action c (or only for agent 1 with EF and SEF). Modifying
 309 the other rewards is however much more costly (under the norm cost function) since each one of
 310 them has L copies of themselves, which requires the same modification by symmetry.

311 7 Conclusion

312 We studied the fairness issue in policy teaching and adopted the notion of envy-freeness to formalize
 313 the problem. Several fundamental questions regarding the existence of a fair solution, the computation
 314 of cost-minimization solution, and the price of considering fairness have been answered in the paper.
 315 For future work, it would be interesting to generalize the model to other reward design settings, where
 316 a larger set of design objectives can be considered. In the case of non-negative adjustment schemes,
 317 we prove that a fair solution exists for one setting. An interesting future research direction would be
 318 to prove this result for a broader class of settings.

319 **Limitations** As we mentioned earlier in the paper, policy teaching is equivalent to reward poisoning
 320 from a technical point of view. Hence, almost any techniques that applies to policy teaching also
 321 applies immediately to solve reward poisoning problems. We note this potential negative social
 322 impact of our results but also remark that since our consideration is fairness we are not aware of any
 323 scenario where a malicious party intends to launch poisoning attack while considering fairness at
 324 the same time. There are many other notions of fairness, equity, and equality. The EF notions we
 325 studied are concerned with the additional rewards provided by the adjustment scheme but not with
 326 the overall rewards. Hence, they are not applicable if the latter should be the key consideration.

327 References

- 328 [1] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press,
329 2018.
- 330 [2] Chen Tessler, Daniel J Mankowitz, and Shie Mannor. Reward constrained policy optimization.
331 In *ICLR*, 2018.
- 332 [3] Deirdre K. Mulligan, Joshua A. Kroll, Nitin Kohli, and Richmond Y. Wong. This thing called
333 fairness: Disciplinary confusion realizing a value in technology. *Proc. ACM Hum.-Comput.*
334 *Interact.*, 3(CSCW), nov 2019.
- 335 [4] Ariel D. Procaccia. Cake cutting: Not just child’s play. *Commun. ACM*, 56(7):78–87, jul 2013.
- 336 [5] Haoqi Zhang and David C. Parkes. Value-based policy teaching with active indirect elicitation.
337 In *AAAI*, pages 208–214, 2008.
- 338 [6] Haoqi Zhang, David C Parkes, and Yiling Chen. Policy teaching through reward function
339 learning. In *EC*, pages 295–304, 2009.
- 340 [7] Kiarash Banihashem, Adish Singla, Jiarui Gan, and Goran Radanovic. Admissible policy
341 teaching through reward design. *arXiv preprint arXiv:2201.02185*, 2022.
- 342 [8] Yuzhe Ma, Xuezhou Zhang, Wen Sun, and Jerry Zhu. Policy poisoning in batch reinforcement
343 learning and control. In *NeurIPS*, pages 14543–14553, 2019.
- 344 [9] Yunhan Huang and Quanyan Zhu. Deceptive reinforcement learning under adversarial manipu-
345 lations on cost signals. In *GameSec*, pages 217–237, 2019.
- 346 [10] Amin Rakhsha, Goran Radanovic, Rati Devidze, Xiaojin Zhu, and Adish Singla. Policy teaching
347 in reinforcement learning via environment poisoning attacks. *CoRR*, abs/2011.10824, 2020.
- 348 [11] Yanchao Sun, Da Huo, and Furong Huang. Vulnerability-aware poisoning mechanism for online
349 rl with unknown dynamics. In *ICLR*, 2021.
- 350 [12] Xuezhou Zhang, Yuzhe Ma, Adish Singla, and Xiaojin Zhu. Adaptive reward-poisoning attacks
351 against reinforcement learning. In *ICML*, pages 11225–11234, 2020.
- 352 [13] Maja J Mataric. Reward functions for accelerated learning. In *ICML*, pages 181–189, 1994.
- 353 [14] Marco Dorigo and Marco Colombetti. Robot shaping: Developing autonomous agents through
354 learning. *Artificial intelligence*, 71(2):321–370, 1994.
- 355 [15] Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transforma-
356 tions: Theory and application to reward shaping. In *ICML*, pages 278–287, 1999.
- 357 [16] Duncan Karl Foley. *Resource allocation and the public sector*. Yale University, 1966.
- 358 [17] G Gamow and M Stern. *Puzzle-math*, edn. *Viking*, 1958.
- 359 [18] Hervé Moulin. *Fair division and collective welfare*. MIT press, 2004.
- 360 [19] Kenneth J Arrow, Amartya Sen, and Kotaro Suzumura. *Handbook of social choice and welfare*,
361 volume 2. Elsevier, 2010.
- 362 [20] Tim Roughgarden. Algorithmic game theory. *Communications of the ACM*, 53(7):78–86, 2010.
- 363 [21] Daniel Halpern and Nisarg Shah. Fair division with subsidy. In *Proceedings of 12th International*
364 *Symposium (SAGT ’19)*, page 374–389, Berlin, Heidelberg, 2019. Springer-Verlag.
- 365 [22] Xiaojin Zhu, Ji Liu, and Manuel Lopes. No learner left behind: On the complexity of teaching
366 multiple learners simultaneously. In *Proceedings of the 26th International Joint Conference on*
367 *Artificial Intelligence, IJCAI ’17*, pages 3588–3594, 2017.
- 368 [23] Teresa Yeo, Parameswaran Kamalaruban, Adish Singla, Arpit Merchant, Thibault Asselborn,
369 Louis Faucon, Pierre Dillenbourg, and Volkan Cevher. Iterative classroom teaching. In
370 *Proceedings of the 33rd AAAI Conference on Artificial Intelligence (AAAI ’19)*, 2019.

- 371 [24] Ritesh Noothigattu, Tom Yan, and Ariel D. Procaccia. Inverse reinforcement learning from
372 like-minded teachers. volume 35, pages 9197–9204, May 2021.
- 373 [25] Paul Duetting, Tim Roughgarden, and Inbal Talgam-Cohen. The complexity of contracts. *SIAM*
374 *Journal on Computing*, 50(1):211–254, 2021.

375 **Checklist**

- 376 1. For all authors...
- 377 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
378 contributions and scope? [Yes]
- 379 (b) Did you describe the limitations of your work? [Yes]
- 380 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 381 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
382 them? [Yes]
- 383 2. If you are including theoretical results...
- 384 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 385 (b) Did you include complete proofs of all theoretical results? [Yes]
- 386 3. If you ran experiments...
- 387 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
388 mental results (either in the supplemental material or as a URL)? [N/A]
- 389 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
390 were chosen)? [N/A]
- 391 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
392 ments multiple times)? [N/A]
- 393 (d) Did you include the total amount of compute and the type of resources used (e.g., type
394 of GPUs, internal cluster, or cloud provider)? [N/A]
- 395 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 396 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 397 (b) Did you mention the license of the assets? [N/A]
- 398 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 399
- 400 (d) Did you discuss whether and how consent was obtained from people whose data you're
401 using/curating? [N/A]
- 402 (e) Did you discuss whether the data you are using/curating contains personally identifiable
403 information or offensive content? [N/A]
- 404 5. If you used crowdsourcing or conducted research with human subjects...
- 405 (a) Did you include the full text of instructions given to participants and screenshots, if
406 applicable? [N/A]
- 407 (b) Did you describe any potential participant risks, with links to Institutional Review
408 Board (IRB) approvals, if applicable? [N/A]
- 409 (c) Did you include the estimated hourly wage paid to participants and the total amount
410 spent on participant compensation? [N/A]