

On Welfare-Centric Fair Reinforcement Learning



Cyrus Cousins



Elita Lobo



Kavosh Asadi



Michael Littman



University of
Massachusetts
Amherst



www.cyruscousins.online/projects/rlfairness/

Reject Egocentrism

Egocentric Viewpoint



- ▶ **A** acts in \mathbb{E} , and \mathbb{E} responds
- ▶ Scalar reward $R(s, a)$ is *intrinsic* to **A**
- ▶ Rational agents selfishly optimize value

$$\operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{\pi, s} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) \mid s_0 \right]$$

Altruistic Viewpoint

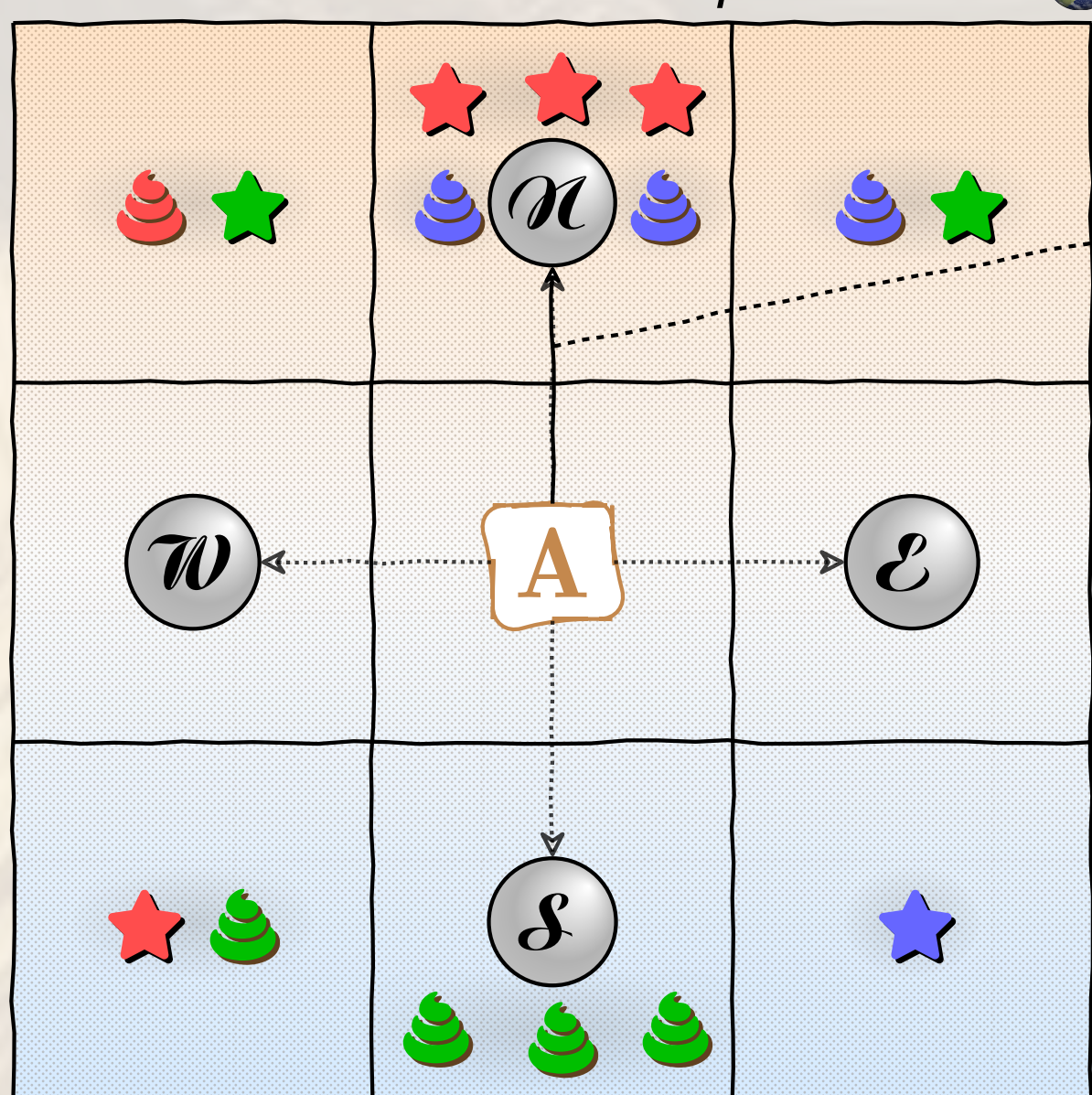


- ▶ **A**'s actions in \mathbb{E} impact *beneficiaries*
- ▶ Vector reward $\mathbf{R}(s, a)$ quantifies impact
- ▶ Altruistic agents optimize *societal welfare*

$$\operatorname{argmax}_{\pi \in \Pi} W \left(i \mapsto \mathbb{E}_{\pi, s} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{R}_i(s_t, \pi(s_t)) \mid s_0 \right] \right)$$

What is Group-Fair Reinforcement Learning?

- ▶ Agent **A** in world \mathbb{E} receives *vector-valued* reward $\mathbf{R}(s, a) \in \mathbb{R}^g$ for g beneficiaries
- ▶ *Beneficiaries* represent *impacted parties*: Individuals, entities, groups, etc.
- ▶ Reward encodes *their response* to **A**- \mathbb{E} interactions



- Yes! +3
- No! -2
- Maybe! +0

Optimize not the value of what **A** wants, but the *welfare* of value functions

$$\operatorname{argmax}_{\pi \in \Pi} W \left(i \mapsto \mathbb{E}_{\pi, s} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{R}_i(s_t, \pi(s_t)) \mid s_0 \right] \right)$$

Geometrically discounted reward

What is a Welfare Function?

- ▶ Utility (value) vector $\mathbf{v} \in \mathbb{R}_{0+}^g$:
- ▶ $W(\mathbf{v}) : \mathbb{R}_{0+}^g \rightarrow \mathbb{R}_{0+}$ *aggregates* utility across beneficiaries

Utilitarian: $W_1(\mathbf{v}) \doteq \frac{1}{g} \sum_{i=1}^g v_i$

Egalitarian: $W_{-\infty}(\mathbf{v}) \doteq \min_{i \in \{1, \dots, g\}} v_i$

Power-Mean: $W_p(\mathbf{v}) \doteq \sqrt[p]{\frac{1}{g} \sum_{i=1}^g v_i^p}$



Even Bandits are Tricky!

"Compromise" 3-Armed Bandit

$\mathbf{R}(s_1, a_1) = \langle 1, 0 \rangle$ $\mathbf{R}(s_1, a_2) = \langle 0, 1 \rangle$

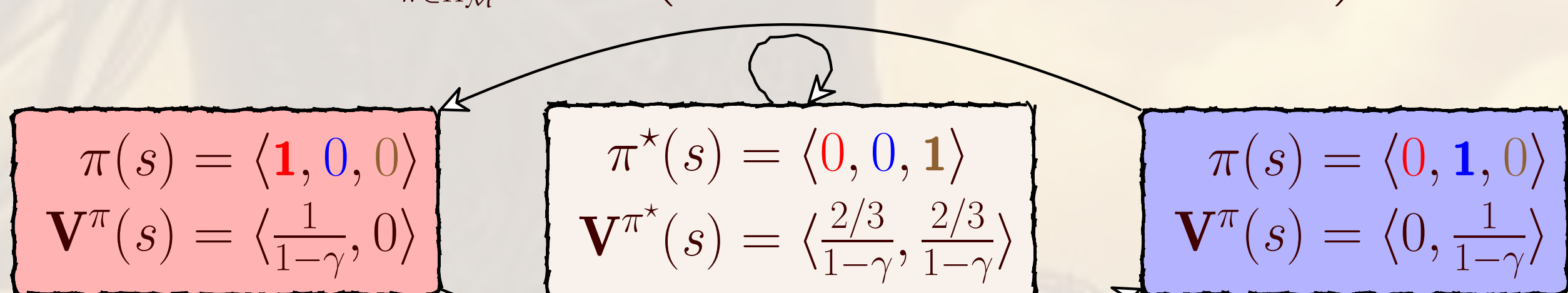
$\pi^1 = \langle \mathbf{1}, 0, 0 \rangle$
 $\pi^2 = \langle 0, \mathbf{1}, 0 \rangle$
 $\pi^* = \langle 0, 0, \mathbf{1} \rangle$

$\mathbf{R}(s_1, a_3) = \langle \frac{2}{3}, \frac{2}{3} \rangle$

Beneficiary policies π^1 and π^2 and fair policy π^* are disjoint!

If $\gamma \geq \frac{1}{2}$: Egalitarian policy iteration oscillates indefinitely

$$\pi^{(t+1)} \leftarrow \operatorname{argmax}_{\pi \in \Pi_M} W_{-\infty} \left(i \mapsto \mathbb{E}_{\pi, s_1} \left[\mathbf{R}_i(s_0, \pi(s_0)) + \gamma \mathbf{V}_i^{\pi^{(t)}}(s_1) \right] \right)$$



On Planning

Policy Iteration

- ✗ Nonconvergent; can oscillate indefinitely

Value Iteration

- ▶ With what Bellman operator? Many obstacles here:
 - ✗ Beneficiaries each have their own *value function* $\mathbf{V}_{1:g}$, but not their own *policy* π
 - ✗ No greedy-optimal substructure (start-state dependence)
- ▶ Planning with *geometrically-discounted state-action occupancy frequencies*

$$\mathbf{d}^* = \operatorname{argmax}_{\mathbf{d} \in \mathbb{R}_{0+}^{S \times A}} W \left(\sum_{s \in S, a \in A} d_{s,a} \mathbf{R}_1(s, a), \sum_{s \in S, a \in A} d_{s,a} \mathbf{R}_2(s, a), \dots, \sum_{s \in S, a \in A} d_{s,a} \mathbf{R}_g(s, a) \right)$$

such that $\forall s \in S : \sum_{a \in A} d_{s,a} = \mathbf{p}_s + \gamma \sum_{s' \in S, a' \in A} \mathbf{P}_s(s', a') d_{s',a'}$,

Take $\pi^*(s, a) \propto d_{s,a}^*$ for all $s \in S, a \in A$

Regret and Mistakes

- ▶ Optimal policy is *stochastic*, can't assess individual actions
- ▶ Assess regret of welfare of *agent policies* $\hat{\pi}_1, \dots, \hat{\pi}_T$

$$\operatorname{Regret}(T) = \sum_{t=1}^T \left(W(\mathbf{V}^{\hat{\pi}_t^*}(s_t)) - W(\mathbf{V}^{\hat{\pi}_t}(s_t)) \right)$$

- ▶ When should we evaluate the agent?
 - ✗ Incoherent to take $s_{t+1} \sim \hat{\pi}_t(s_t)$
 - ▶ *Geometric discounting* suggests *geometric episode length*
 - ▶ Unfair to execute each $\hat{\pi}_t(s_t)$ (start-state dependence)
 - ? **Continuous:** Follow $\hat{\pi}_t$ for Geometric $(1 - \gamma)$ steps, resume
 - ? **Episodic:** End episode, draw s_{t+1} from start-state distribution
- ▶ A policy $\hat{\pi}$ is a *mistake* at s if $W(\mathbf{V}^{\hat{\pi}^*}(s)) - W(\mathbf{V}^{\hat{\pi}}(s)) > \epsilon$
 - ✗ Exploration actions are probably mistakes
 - ? Can exploitation confidently avoid mistakes?



Learning Model: KWIK-AF

The KWIK-AF Learner

Know-What-It-Knows

Adversarial-Fair

- ▶ **KWIK Learner:** At each step, agent has two choices:
 - Output an ϵ -optimal *exploitation policy* π_{xpt}
 - ✗ With probability at least $1 - \delta$, for *all time*
 - ✗ *No mistakes:* $W(\mathbf{V}^{\hat{\pi}^*}(s)) - W(\mathbf{V}^{\pi_{\text{xpt}}}(s)) > \epsilon$
 - Output an *exploration action* a
 - ✓ Receive (s, a, \mathbf{r}, s') tuple, return control to agent in s'
 - ✗ Limited budget: Only $m(|S|, |A|, \gamma, R_{\max}, g, \epsilon, \delta)$ exploration actions, *ever*
- ▶ **Adversarial-Fair:** Algorithm must be *flexible* and *robust*
 - ▶ Optimize for *adversarially selected* welfare function $W(\cdot)$ at each step
 - ▶ When **A** outputs a policy π_{xpt} :
 - ▶ Move **A** to *adversarial* s' , provide *no feedback!*

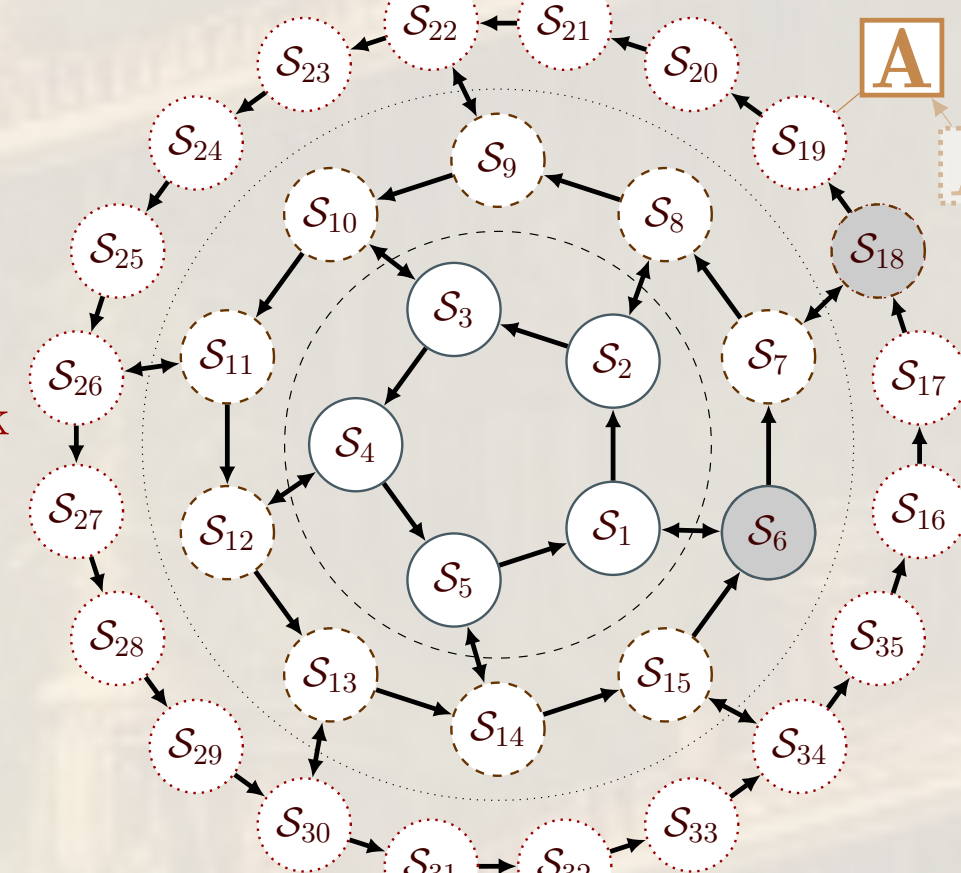
Don't make a mistake.

You may ask a few questions

— but you must learn KWIK.

E⁴: The Equitable Explicit Explore Exploit Algorithm

- ▶ Partition state space into three sets: $\mathcal{S}_{\text{unk}}, \mathcal{S}_{\text{out}}, \mathcal{S}_{\text{inn}}$
- ▶ **Unknown** \mathcal{S}_{unk} : Insufficient samples (fewer than m_{knw}) to estimate reward $\mathbf{R}(s, a)$ and transition $\mathbf{P}(s)$
- ▶ **Outer-Known** \mathcal{S}_{out} : Some escape policy π_{esc} can reach \mathcal{S}_{unk} in T steps with probability at least E
- ▶ **Inner-Known** \mathcal{S}_{inn} : No policy can reach \mathcal{S}_{unk} in T steps with probability at least E
- ▶ Learning moves states from $\mathcal{S}_{\text{unk}} \rightarrow \mathcal{S}_{\text{out}} \rightarrow \mathcal{S}_{\text{inn}}$



~ The E⁴ Algorithm ~

- If in \mathcal{S}_{unk} : Explore, observe (s, a, \mathbf{r}, s') , update empirical MDP \hat{M} , update $\mathcal{S}_{\text{unk}}, \mathcal{S}_{\text{out}}, \mathcal{S}_{\text{inn}}$
- If escape in progress: Follow π_{esc} and decrement timer
- If in \mathcal{S}_{out} : Begin T -step escape attempt in $\pi_{\text{esc}} \leftarrow \operatorname{argmax}_{\pi \in \Pi_T} \sum_{s \in \mathcal{S}} \sum_{s_{t+1} \sim \mathbf{P}(s_t, \pi(s_t, t))} \mathbb{P} \left(\bigvee_{i=0}^T s_i \in \mathcal{S}_{\text{unk}} \mid s_0 = s \right)$
- Otherwise in \mathcal{S}_{inn} : Output exploit policy $\pi_{\text{xpt}} \leftarrow \operatorname{argmax}_{\pi \in \Pi_M} W(\hat{\mathbf{V}}^{\pi}(s))$

E⁴ Theory

- ▶ Can set T, E, m_{knw} to ϵ - δ KWIK-AF learn
 - ▶ At any point in the execution of E⁴, **A** can act effectively:
 - Can exploit from \mathcal{S}_{inn}
 - Can explore directly from \mathcal{S}_{unk}
 - Can explore indirectly from \mathcal{S}_{out}
 - ▶ Escape succeeds with some probability
 - ▶ E⁴ KWIK-AF learns w.r.t. the class of all $\lambda \|\cdot\|_{\infty}$ Lipschitz-continuous welfare functions
- Exploration Budget: $m(|S|, |A|, \gamma, R_{\max}, g, \epsilon, \delta) \in \mathbf{O} \left(|S|^2 |A| \left(\frac{\lambda R_{\max}}{\epsilon(1-\gamma)} \log \frac{\lambda R_{\max}}{\epsilon(1-\gamma)} \right)^3 \log \frac{|S||A|g}{\delta} \right)$
 $\subseteq \text{Poly} \left(|S|, |A|, \frac{1}{1-\gamma}, R_{\max}, \log g, \frac{1}{\epsilon}, \log \frac{1}{\delta}, \lambda \right)$

In Summary

- ▶ From Egocentric to Altruistic Agents
 - ▶ Agent **A** acts in \mathbb{E} , impacting *beneficiaries*
 - ▶ Vector-valued (per-beneficiary) reward $\mathbf{R}(s, a)$
 - ▶ Social planner's problem:
 - ▶ Optimize *welfare of value functions* $\operatorname{argmax}_{\pi \in \Pi} W(\mathbf{V}^{\pi}(s))$
 - ▶ Incorporate *fairness* into *sequential learning* problems
- ▶ KWIK-AF: A Model of Fair RL
 - ▶ Adversarial flexibility
 - ▶ Societal welfare objectives
 - ▶ Tolerate *no mistakes*, allow *bounded exploration*
 - ▶ Challenging model of learning, subsumes PAC-MDP
- ▶ Efficient Learning and Planning
 - ▶ Learn with E⁴: Poly(\dots) exploration budget
 - ▶ Plan with *convex programming* on *state-action measure*
 - ▶ Fair RL and classic RL are comparably difficult

