

On Welfare-Centric Fair Reinforcement Learning



Cyrus Cousins



Elita Lobo



Kavosh Asadi



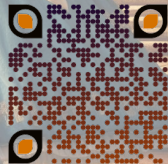
Michael Littman



Reinforcement
Learning
Conference



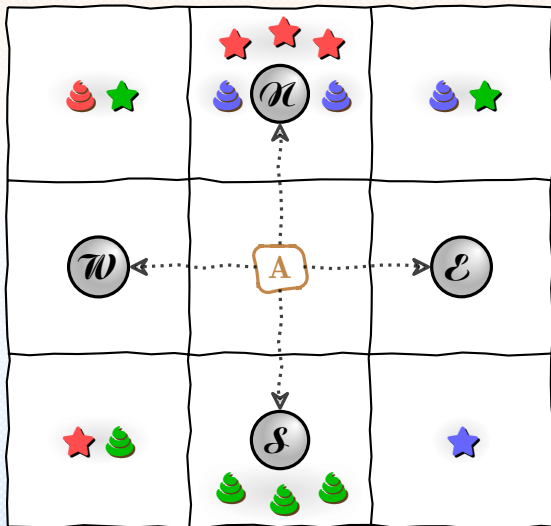
University of
Massachusetts
Amherst



www.cyruscousins.online/projects/rlfairness/

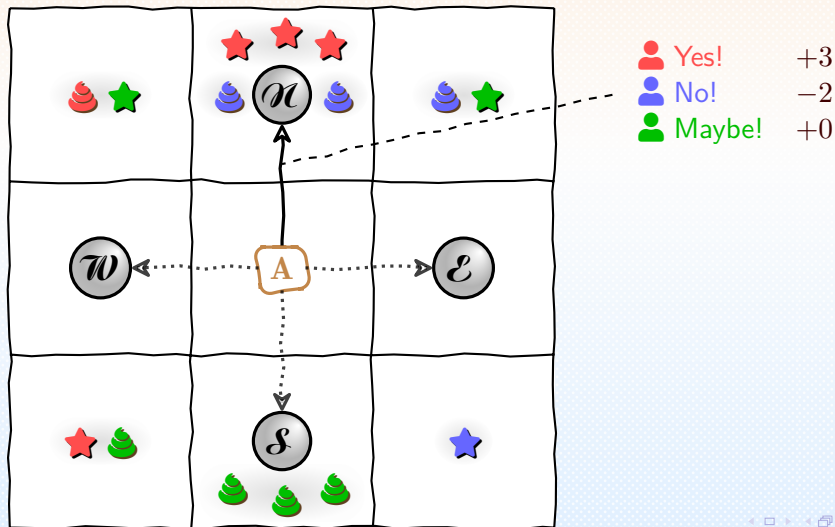
What is Group-Fair Reinforcement Learning?

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





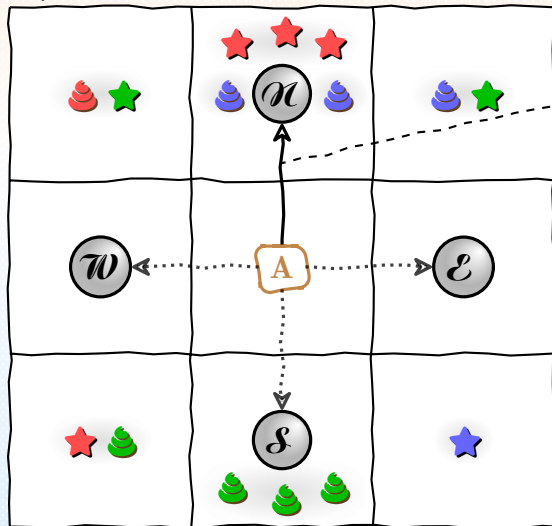
What is Group-Fair Reinforcement Learning?




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



What is Group-Fair Reinforcement Learning?

- ▶ Agent **A** in world  receives *vector-valued* reward $\mathbf{R}(s, a) \in \mathbb{R}^g$ from g beneficiaries
 - ▶ *Beneficiaries* represent *impacted parties*: Individuals, entities, groups, etc.
 - ▶ Reward encodes *their response* to **A**- interactions
- ▶ Optimize not the value of *what A wants*, but the *welfare* of beneficiary value functions





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

Objective:

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

Egocentric Viewpoint





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

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

Reject Egocentrism

Egocentric Viewpoint




- ▶ **A** acts in , and  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  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 a Welfare Function?

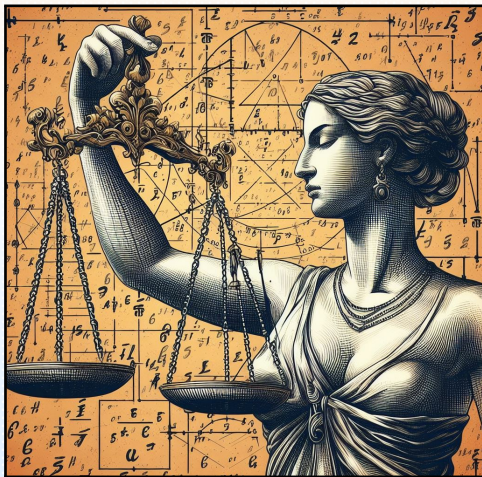
- ▶ Given g beneficiaries
- ▶ Utility (value) vector $v \in \mathbb{R}_{0+}^g$

$$v = \langle \text{★★★}, \text{★★}, \text{★} \rangle$$



What is a Welfare Function?

- ▶ Given g beneficiaries
- ▶ Utility (value) vector $v \in \mathbb{R}_{0+}^g$
$$v = \left\langle \text{★★★}, \text{★★}, \text{★} \right\rangle$$
- ▶ $W(v) : \mathbb{R}_{0+}^g \rightarrow \mathbb{R}_{0+}$ aggregates utility across beneficiaries
- ▶ *Welfare functions encode social values*

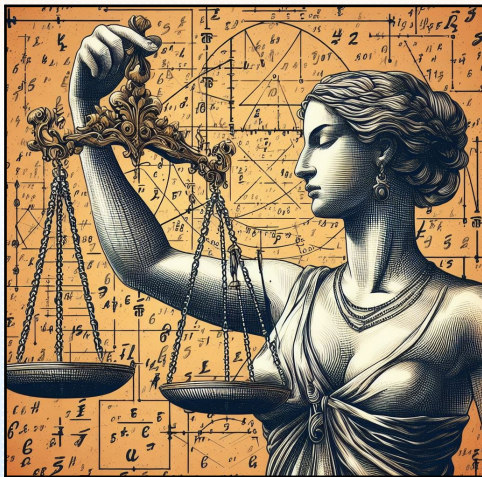


What is a Welfare Function?

- ▶ Given g beneficiaries
- ▶ Utility (value) vector $\mathbf{v} \in \mathbb{R}_{0+}^g$
$$\mathbf{v} = \left\langle \text{★★★}, \text{★★}, \text{★} \right\rangle$$
- ▶ $W(\mathbf{v}) : \mathbb{R}_{0+}^g \rightarrow \mathbb{R}_{0+}$ aggregates utility across beneficiaries
- ▶ *Welfare functions encode social values*

- ▶ Common welfare functions

- ▶ Utilitarian: $W_1(\mathbf{v}) \doteq \frac{1}{g} \sum_{i=1}^g v_i$
- ▶ Egalitarian: $W_{-\infty}(\mathbf{v}) \doteq \min_{i \in \{1, \dots, g\}} v_i$
- ▶ p Power-Mean: $W_p(\mathbf{v}) \doteq \sqrt[p]{\frac{1}{g} \sum_{i=1}^g v_i^p}$

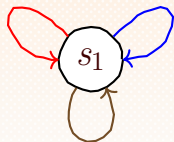


Even Bandits are Tricky!

“Compromise” 3-Armed Bandit

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

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

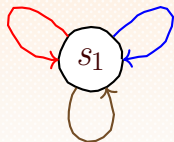


$$\mathbf{R}(s_1, a_3) = \langle \frac{2}{3}, \frac{2}{3} \rangle$$

Even Bandits are Tricky!

“Compromise” 3-Armed Bandit

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



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

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

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

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

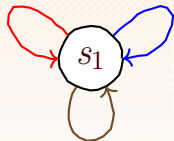
$$\pi^* = \langle 0, 0, \mathbf{1} \rangle$$

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

Even Bandits are Tricky!

“Compromise” 3-Armed Bandit

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



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

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

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

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

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

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

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

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

$$\begin{aligned} \pi(s) &= \langle \mathbf{1}, \mathbf{0}, \mathbf{0} \rangle \\ \mathbf{V}^{\pi}(s) &= \langle \frac{1}{1-\gamma}, \mathbf{0} \rangle \end{aligned}$$

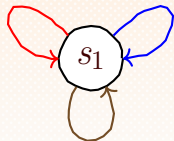
$$\begin{aligned} \pi(s) &= \langle \mathbf{0}, \mathbf{1}, \mathbf{0} \rangle \\ \mathbf{V}^{\pi}(s) &= \langle \mathbf{0}, \frac{1}{1-\gamma} \rangle \end{aligned}$$

Even Bandits are Tricky!

“Compromise” 3-Armed Bandit

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

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



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

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

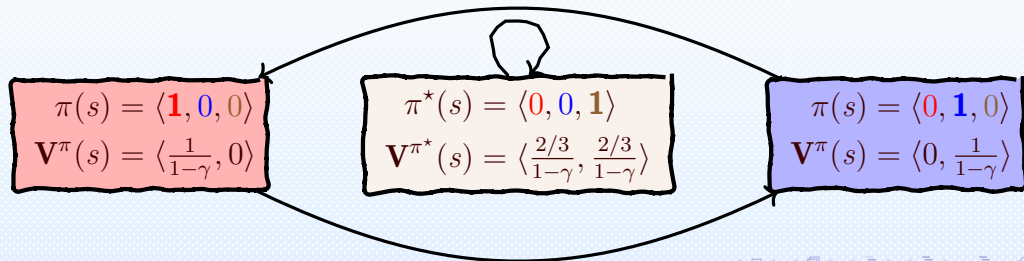
$$\pi^2 = \langle 0, \mathbf{1}, 0 \rangle$$

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

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

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

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

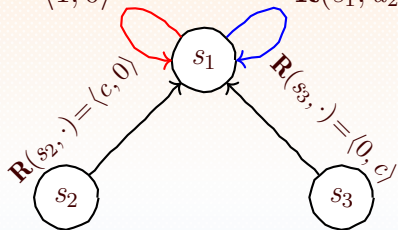


Overcoming Initial Disparity

“Asymmetric Start Bandit” MDP

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

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

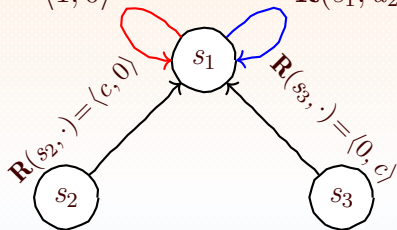


Overcoming Initial Disparity

“Asymmetric Start Bandit” MDP

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

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



$$\pi^1(s_1) = \langle \mathbf{1}, \mathbf{0} \rangle$$

$$\pi^2(s_1) = \langle \mathbf{0}, \mathbf{1} \rangle$$

$$\pi^*(s_1 \text{ from } s_2) = \langle \frac{1}{2} - \frac{1-\gamma}{2\gamma} c, \frac{1}{2} + \frac{1-\gamma}{2\gamma} c \rangle$$

$$\pi^*(s_1 \text{ from } s_3) = \langle \frac{1}{2} + \frac{1-\gamma}{2\gamma} c, \frac{1}{2} - \frac{1-\gamma}{2\gamma} c \rangle$$

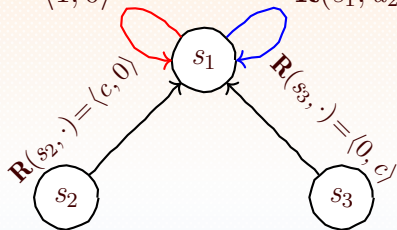
Fair policy π^* is *start-state dependent!*

Overcoming Initial Disparity

“Asymmetric Start Bandit” MDP

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

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



$$\pi^1(s_1) = \langle \mathbf{1}, \mathbf{0} \rangle$$

$$\pi^2(s_1) = \langle \mathbf{0}, \mathbf{1} \rangle$$

$$\pi^*(s_1 \text{ from } s_2) = \langle \frac{1}{2} - \frac{1-\gamma}{2\gamma} c, \frac{1}{2} + \frac{1-\gamma}{2\gamma} c \rangle$$

$$\pi^*(s_1 \text{ from } s_3) = \langle \frac{1}{2} + \frac{1-\gamma}{2\gamma} c, \frac{1}{2} - \frac{1-\gamma}{2\gamma} c \rangle$$

Fair policy π^* is *start-state dependent!*

Lemma (Optimality of Stationary Policies)

For any start state $s_0 \in \mathcal{S}$, there exists some $W(\cdot)$ -optimal policy

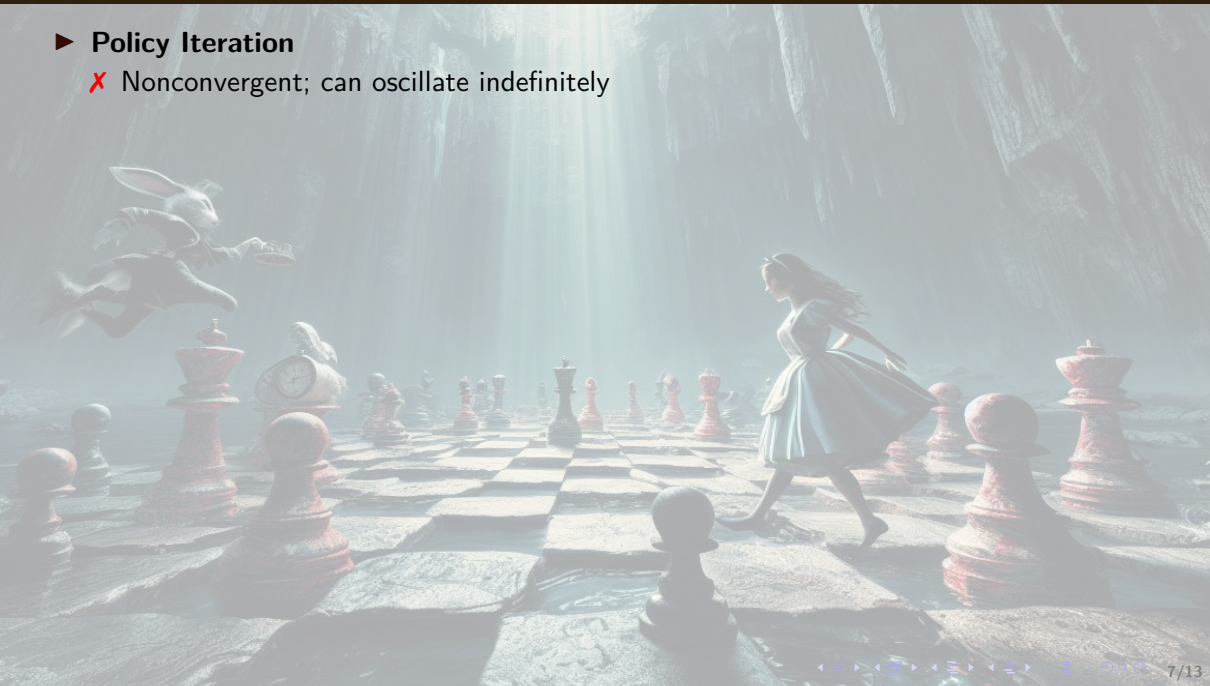
$$\pi^* \in \operatorname{argmax}_{\pi \in \Pi_{\mathcal{M}}} W \left(\mathbf{V}_1^\pi(s_0), \dots, \mathbf{V}_g^\pi(s_0) \right)$$

that is a stationary (Markovian) stochastic policy

On Planning

► Policy Iteration

✗ Nonconvergent; can oscillate indefinitely



On Planning

► Policy Iteration

✗ Nonconvergent; can oscillate indefinitely

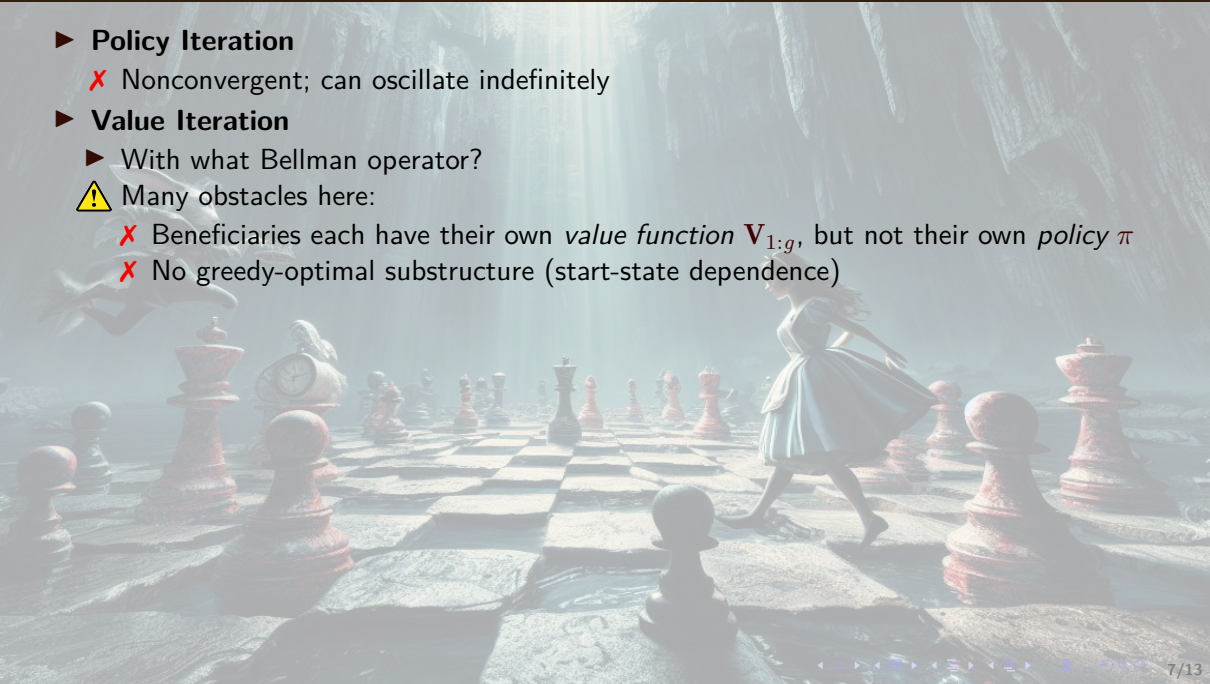
► Value Iteration

► With what Bellman operator?

⚠ Many obstacles here:

✗ Beneficiaries each have their own *value function* $V_{1:g}$, but not their own *policy* π

✗ No greedy-optimal substructure (start-state dependence)



On Planning

► Policy Iteration

✗ Nonconvergent; can oscillate indefinitely

► Value Iteration

► With what Bellman operator?

⚠ Many obstacles here:

✗ Beneficiaries each have their own *value function* $V_{1:g}$, but not their own *policy* π

✗ No greedy-optimal substructure (start-state dependence)

💡 Planning with *geometrically-discounted state-action occupancy frequencies*

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

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

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

✓ Approximately optimize π^* with *convex programming*

Regret and Mistakes

- ▶ Optimal policy is *stochastic*, can't assess individual actions



Assess regret of welfare of *agent policies* $\hat{\pi}_1, \dots, \hat{\pi}_T$

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



Regret and Mistakes

- ▶ Optimal policy is *stochastic*, can't assess individual actions



Assess regret of welfare of *agent policies* $\hat{\pi}_1, \dots, \hat{\pi}_T$

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

- ▶ 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)

? **Episodic:** End episode, draw s_{t+1} from start-state distribution



Regret and Mistakes

- ▶ Optimal policy is *stochastic*, can't assess individual actions

💡 Assess regret of welfare of *agent policies* $\hat{\pi}_1, \dots, \hat{\pi}_T$

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

- ▶ 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)

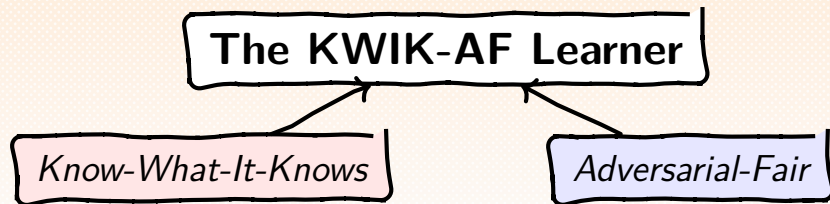
? **Episodic:** End episode, draw s_{t+1} from start-state distribution

- ▶ A policy $\hat{\pi}$ is a *mistake* at s if $W(\mathbf{V}^{\pi^*}(s)) - W(\mathbf{V}^{\hat{\pi}}(s)) > \varepsilon$

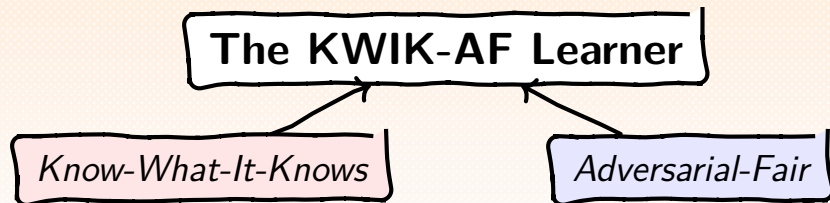
✗ Exploration actions are probably mistakes

? Can exploitation confidently avoid mistakes?





- ▶ **KWIK Learner:** At each step, in state s , \mathbf{A} can either
 1. Output an ε -optimal *exploitation policy* π_{xpt}
 - ✗ With probability at least $1 - \delta$, for all time
 - ✗ No mistakes: $W(\mathbf{V}^{\pi_s^*}(s)) - W(\mathbf{V}^{\pi_{\text{xpt}}}(s)) > \varepsilon$
 2. Output an *exploration action* a
 - ✓ Receive (s, a, r, s') tuple, return control to agent in s'
 - ✗ Limited budget: Only $m(|\mathcal{S}|, |\mathcal{A}|, \gamma, R_{\max}, g, \varepsilon, \delta)$ exploration actions



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



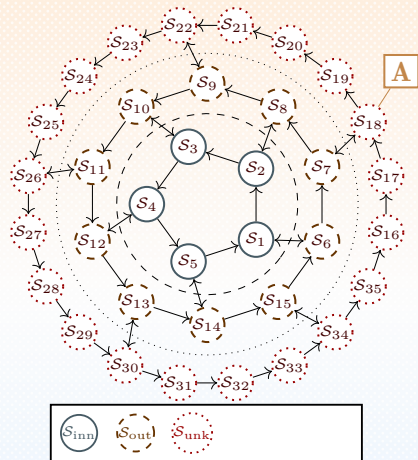
Don't make a mistake.

You may ask a few questions

— but you must learn KWIK.

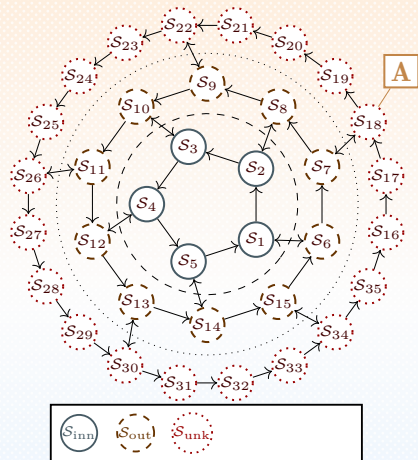
E⁴: The Equitable Explicit Explore Exploit Algorithm

- ▶ Partition state space into three sets: \mathcal{S}_{unk} , \mathcal{S}_{out} , \mathcal{S}_{inn}
- ▶ **Unknown** \mathcal{S}_{unk} : Insufficient samples ($\leq m_{\text{knw}}$) to estimate reward $\mathbf{R}(s, a)$ and transition $\mathbf{P}(s)$



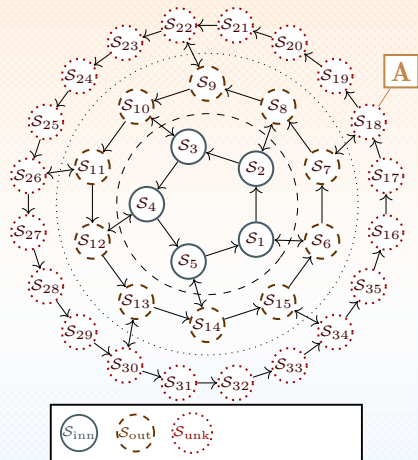
E⁴: The Equitable Explicit Explore Exploit Algorithm

- ▶ Partition state space into three sets: \mathcal{S}_{unk} , \mathcal{S}_{out} , \mathcal{S}_{inn}
- ▶ **Unknown** \mathcal{S}_{unk} : Insufficient samples ($\leq m_{\text{knw}}$) 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



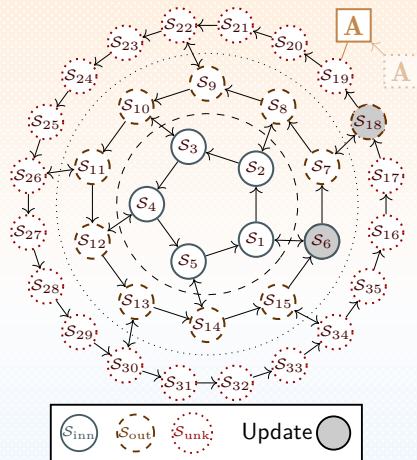
E⁴: The Equitable Explicit Explore Exploit Algorithm

- ▶ Partition state space into three sets: \mathcal{S}_{unk} , \mathcal{S}_{out} , \mathcal{S}_{inn}
- ▶ **Unknown** \mathcal{S}_{unk} : Insufficient samples ($\leq m_{\text{knw}}$) 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



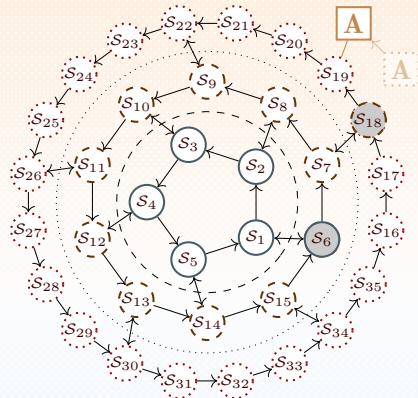
E⁴: The Equitable Explicit Explore Exploit Algorithm

- ▶ Partition state space into three sets: \mathcal{S}_{unk} , \mathcal{S}_{out} , \mathcal{S}_{inn}
- ▶ **Unknown** \mathcal{S}_{unk} : Insufficient samples ($\leq m_{\text{knw}}$) 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 $\mathcal{S}_{\text{unk}} \rightarrow \mathcal{S}_{\text{out}} \rightarrow \mathcal{S}_{\text{inn}}$



E⁴: The Equitable Explicit Explore Exploit Algorithm

- ▶ Partition state space into three sets: \mathcal{S}_{unk} , \mathcal{S}_{out} , \mathcal{S}_{inn}
- ▶ **Unknown** \mathcal{S}_{unk} : Insufficient samples ($\leq m_{\text{knw}}$) 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 ~

1. If in \mathcal{S}_{unk} : Explore, observe (s, a, r, s') , update empirical MDP $\hat{\mathcal{M}}$, update \mathcal{S}_{unk} , \mathcal{S}_{out} , \mathcal{S}_{inn}
2. If in \mathcal{S}_{out} : Begin escape attempt (follow π_{esc} for T steps)

$$\pi_{\text{esc}} \leftarrow \operatorname{argmax}_{\pi \in \Pi_T} \sum_{s \in \mathcal{S}} \mathbb{P}(s_{t+1} \sim \hat{\mathbf{P}}(s_t, \pi(s_t, t))) \left(\bigvee_{i=0}^T s_i \in \mathcal{S}_{\text{unk}} \mid s_0 = s \right)$$

3. Otherwise in \mathcal{S}_{inn} : Output exploit policy $\pi_{\text{xpt}} \leftarrow \operatorname{argmax}_{\pi \in \Pi_{\mathcal{M}}} W(\hat{\mathbf{V}}^\pi(s))$



Lemma (Explore-Exploit)

At any point in the execution of E^4 , \mathbf{A} can act effectively:

1. Can exploit from \mathcal{S}_{inn}
2. Can explore directly from \mathcal{S}_{unk}
3. Can explore indirectly from \mathcal{S}_{out} (escape succeeds with some probability)



Lemma (Explore-Exploit)

At any point in the execution of E^4 , \mathbf{A} can act effectively:

1. Can exploit from \mathcal{S}_{inn}
2. Can explore directly from \mathcal{S}_{unk}
3. Can explore indirectly from \mathcal{S}_{out} (escape succeeds with some probability)

Theorem (E^4 is a KWIK-AF Learner)


E^4 is a KWIK-AF learner w.r.t. the class of all $\lambda\|\cdot\|_\infty$ Lipschitz-continuous welfare functions, with exploration budget

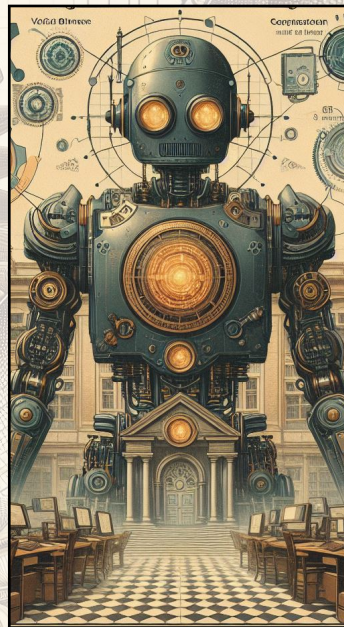
$$m(|\mathcal{S}|, |\mathcal{A}|, \gamma, R_{\max}, g, \varepsilon, \delta) \in \mathbf{O} \left(|\mathcal{S}|^2 |\mathcal{A}| \left(\frac{\lambda R_{\max}}{\varepsilon(1-\gamma)} \log_{\frac{1}{\gamma}} \left(\frac{\lambda R_{\max}}{\varepsilon(1-\gamma)} \right) \right)^3 \log \frac{|\mathcal{S}| |\mathcal{A}| g}{\delta} \right)$$

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

In Summary


▶ From Ego-centric to Altruistic Agents

- ▶ Agent **A** acts in , impacting *beneficiaries*
- ▶ Vector-valued (per-beneficiary) reward $\mathbf{R}(s, a)$
- ▶ Social planner's problem:
 - ▶ Optimize *welfare* of *value functions* $\operatorname{argmax}_{\pi \in \Pi} W(\mathbf{V}^{\pi}(s))$



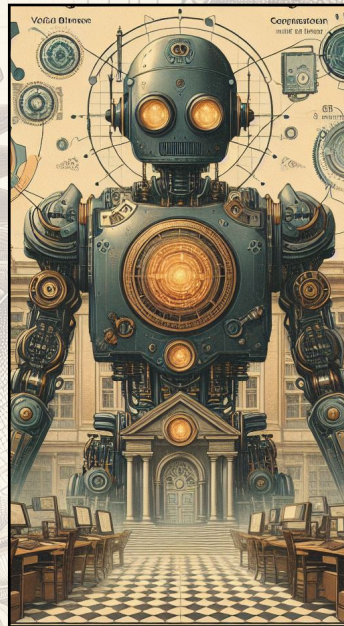
In Summary

▶ From Ego-centric to Altruistic Agents

- ▶ Agent **A** acts in , impacting *beneficiaries*
- ▶ Vector-valued (per-beneficiary) reward $\mathbf{R}(s, a)$
- ▶ Social planner's problem:
 - ▶ Optimize *welfare* of *value functions* $\operatorname{argmax}_{\pi \in \Pi} W(\mathbf{V}^{\pi}(s))$


▶ KWIK-AF: A Model of Fair RL

- ▶ Adversarial flexibility
- ▶ Societal welfare objectives
- ▶ No mistakes from bounded exploration



In Summary

▶ From Ego-centric to Altruistic Agents

- ▶ Agent **A** acts in , impacting *beneficiaries*
- ▶ Vector-valued (per-beneficiary) reward $\mathbf{R}(s, a)$
- ▶ Social planner's problem:
 - ▶ Optimize *welfare* of *value functions* $\operatorname{argmax}_{\pi \in \Pi} W(\mathbf{V}^{\pi}(s))$

▶ KWIK-AF: A Model of Fair RL

- ▶ Adversarial flexibility
- ▶ Societal welfare objectives
- ▶ No mistakes from bounded exploration

▶ Efficient Learning and Planning

- ▶ KWIK-AF learn with E^4
- ▶ Plan with convex programming on state-action measure
- ▶ Polynomial exploration budget, time complexity

