

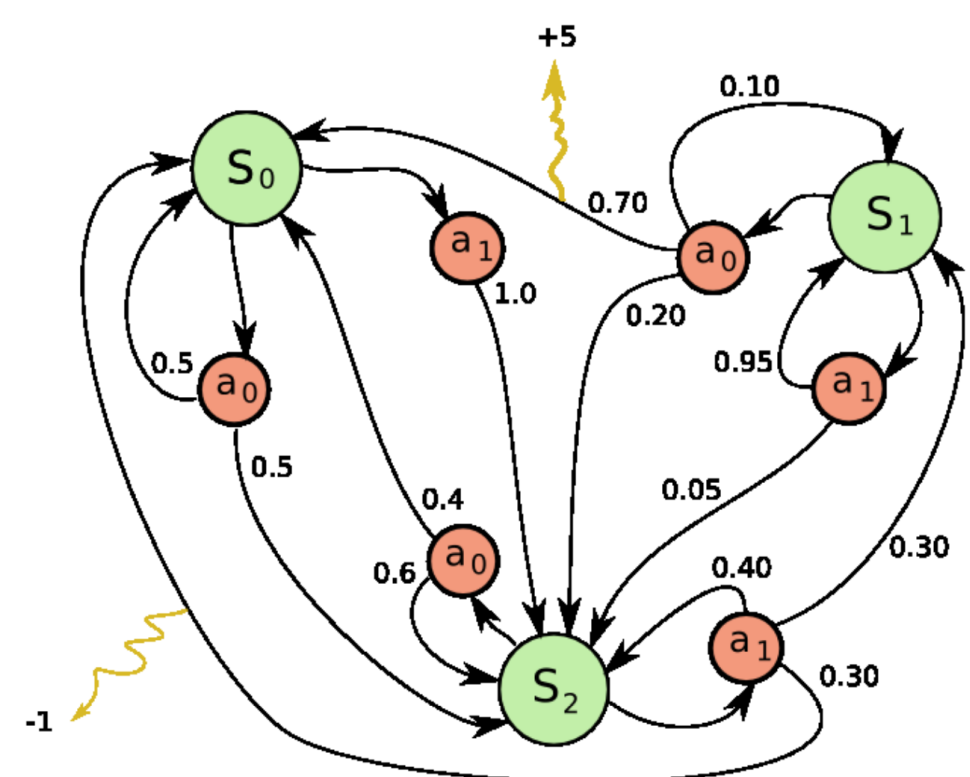
① What type of learning

Environment: MDP, **unknown dynamics**, **unknown cost**

INPUT: A finite set of **expert demonstrations**

GOAL: Learn a policy that performs at least as *good* as the expert

② Markov decision processes (MDPs)



- **Markov decision model** $\mathcal{M}_c \triangleq (\mathcal{X}, \mathcal{A}, P, \gamma, \nu_0, c)$
- $P(x'|x, a) = \text{Prob}(x_{t+1} = x' | x_t = x, a_t = a)$,
- Π_0 set of **stationary Markov policies** π , $\pi(a|x) = \text{Prob}(a_t = a | x_t = x)$,
- $\nu_0 \in \Delta_S$ **initial state distribution**, $c \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ cost function, $\gamma \in (0, 1)$ discount factor.
- $a_t \sim \pi(\cdot | x_t)$; $x_{t+1} \sim P(\cdot | x_t, a_t)$; $c(x_t, a_t)$
- Occupancy measure μ_π induced by a policy

$$\mu_\pi(x, a) := (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(x_t = x, a_t = a | x_0 \sim \nu_0, \pi)$$

- **Minimize a cost criterion**

$$\min_{\pi \in \Pi} \rho_c(\pi), \quad \text{where} \quad \rho_c(\pi) \triangleq (1 - \gamma) \mathbb{E}_{\nu_0}^{\pi} \left[\sum_{t=0}^{\infty} \gamma^t c(x_t, a_t) \right].$$

- **LP formulation**

$$\begin{aligned} \min_{\mu \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}} \max_{c \in \mathcal{C}} \langle \mu - \mu_{\pi_E}, c \rangle \\ \text{s.t. } E^T \mu = (1 - \gamma) \nu_0 + \gamma P^T \mu, \quad \mu \geq 0. \end{aligned} \quad (\text{Primal IL})$$

- **Linear MDP assumption** There exists mappings $\phi : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}^m$ and $g : \mathcal{X} \rightarrow \mathbb{R}^m$ and a vector $w \in \mathcal{W} := \{w \in \mathbb{R}^m : \|w\|_2 \leq 1\}$ such that

$$c(s, a) = \langle \phi(s, a), w \rangle \quad P(s' | s, a) = \langle \phi(s, a), g(s') \rangle$$

that is, in matrix form

$$c = \Phi w \quad P = \Phi M$$

③ The constraint splitting trick

- We plug in the **(Linear MDP)** structure in **(Primal IL)** as follows. A similar trick appeared outside the imitation learning in (Mehta and Meyn, 2020), (Lee and He, 2019) and (Bas-Serrano et al., 2021).

$$\begin{aligned} \min_{\mu \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}} \max_{w \in \mathcal{W}} \langle \mu - \mu_{\pi_E}, \Phi w \rangle \\ \text{s.t. } E^T \mu = (1 - \gamma) \nu_0 + \gamma M^T \Phi^T \mu \Rightarrow \end{aligned} \quad \begin{aligned} \min_{\lambda \in \Delta^m, \mu \in \mathbb{R}^{\mathcal{X} \times \mathcal{A}}} \max_{w \in \mathcal{W}} \langle \lambda - \Phi^T \mu^{\pi_E}, w \rangle \\ \text{s.t. } E^T \mu - \gamma M^T \lambda = (1 - \gamma) \nu_0 \\ \Phi^T \mu = \lambda \end{aligned} \quad (\text{Primal' IL})$$

- Our algorithm consists in applying inexact proximal point updates for μ and λ on the Lagrangian of **Primal' IL**.

④ The algorithm

Proximal Point Imitation Learning: P²IL

Initialize π_0 as uniform distribution over \mathcal{A}

for $k = 1, \dots, K$ **do**

Policy evaluation :

$$(w_k, \theta_k) \approx \text{argmax}_{w \in \mathcal{W}, \theta \in \Theta} \mathcal{G}_k(w, \theta)$$

Policy improvement :

$$\pi_k(a|s) \propto \pi_{k-1}(a|s) e^{-\alpha Q_{\theta_k}(s, a)}$$

Where $\mathcal{G}_k(w, \theta)$ is called (negative) logistic Bellman error (Bas-Serrano et al., 2021) and it is the following concave and smooth function. ^a

$$-\frac{1}{\eta} \log \sum_{i=1}^m (\Phi^T \mu_{k-1})(i) e^{-\eta \delta_{w, \theta}^k(i)} + (1 - \gamma) \langle \nu_0, V_{\theta}^k \rangle - \langle \lambda_{\pi_E}, \Phi^T w \rangle,$$

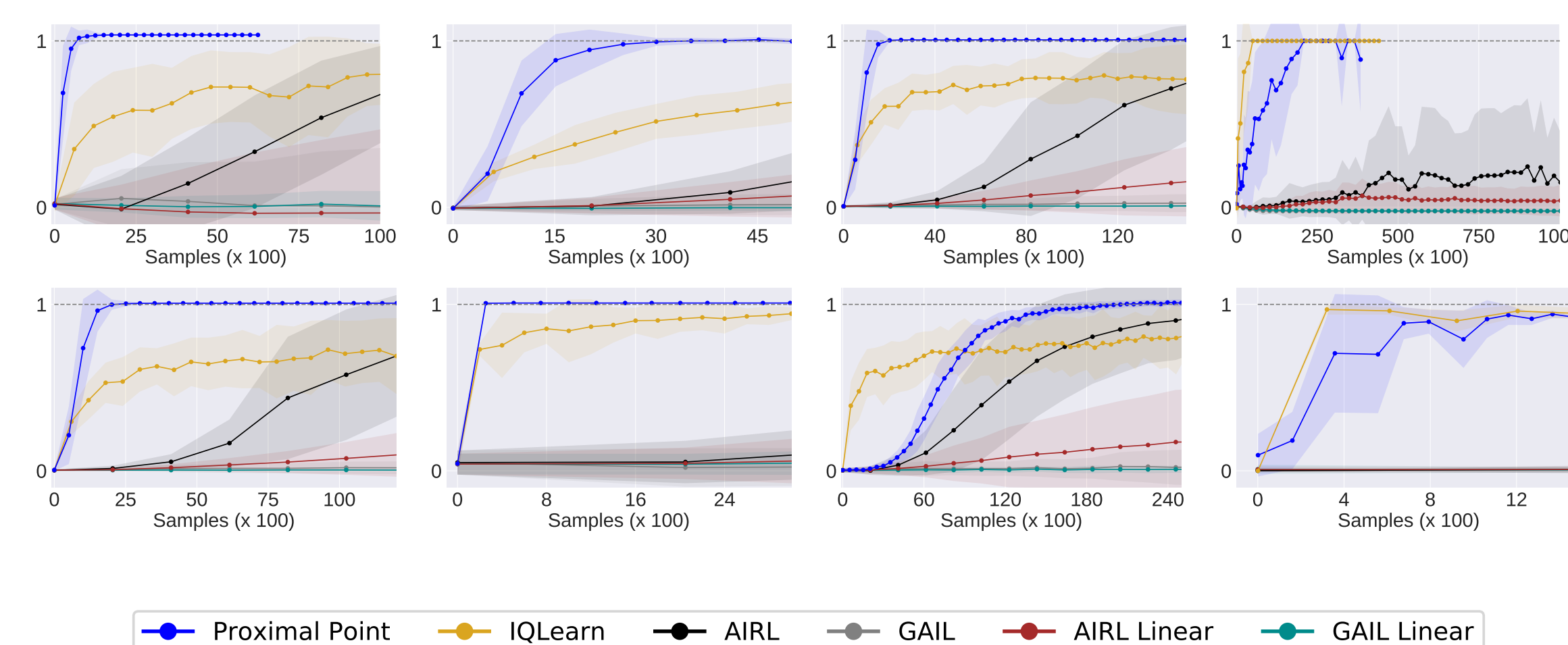
^aWe use $\delta_{w, \theta}^k \triangleq w + \gamma M V_{\theta}^k - \theta$ and $V_{\theta}^k \triangleq -\frac{1}{\alpha} \log (\sum_a \pi_{\mu_{k-1}}(a|s) e^{-\alpha Q_{\theta}(s, a)})$ and $Q_{\theta} = \Phi \theta$.

Theorem 1 (Resources Guarantees) Let us define the \mathcal{C} -distance between π and π_E , $d_{\mathcal{C}}(\pi, \pi_E) \triangleq \max_{c \in \mathcal{C}} (\rho_c(\pi) - \rho_c(\pi_E))$. Using $\Omega(KT) = \Omega(\epsilon^{-5})$ sample transitions, $\Omega(\epsilon^{-2})$ expert trajectories and approximately solving $\Omega(\epsilon^{-1})$ concave maximization problems, we can ensure $d_{\mathcal{C}}(\hat{\pi}, \pi_E) \leq \mathcal{O}(\epsilon + \varepsilon)$, with high probability.

- We consider errors in the maximization of $\mathcal{G}_k(w, \theta)$, i.e. $\epsilon_k = \mathcal{G}_k(w_k^*, \theta_k^*) - \mathcal{G}_k(w_k, \theta_k)$.
- **First**, we show how errors propagate. **Second**, we control that the errors are small using a Biased Stochastic Gradient Ascent subroutine.

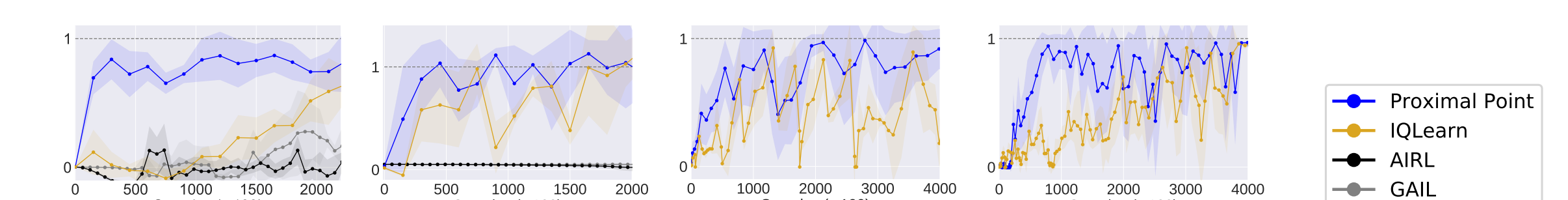
⑤ Discrete Actions Experiments

- From the left to right: WideTree, RiverSwim, SingleChain, DoubleChain, Cartpole, Two State, Gridworld and Acrobot.



⑥ Continuous Control Experiments

- From the left to right: HalfCheetah, Ant, Hopper, Walker.



- Nonlinear function approximation and continuous actions are not covered by our theory.
- However the empirical performance is convincing vs the state-of-the-art IQLearn (Garg et al., 2021).



Full paper