



Problem Formulation

→ Imitation Learning (IL): mimic the expert π^E from demonstrations:

$$\mathcal{D} = \{\text{tr} = (s_1, a_1, s_2, a_2, \dots, s_H, a_H); a_h \sim \pi_E(\cdot|s_h)\}.$$

→ Adversarial imitation learning (AIL): imitate by state-action distribution matching

$$\min_{\pi \in \Pi} \sum_{h=1}^H \left\| d_h^{\pi, P} - \tilde{d}_h^{\pi^E} \right\|_1 \quad (1)$$

- $d_h^{\pi, P}(s, a) = \mathbb{P}^\pi(s_h = s, a_h = a)$: true state-action distribution.
- $\tilde{d}_h^{\pi^E}$: estimation of $d_h^{\pi^E}$ from the finite dataset \mathcal{D} .

→ Setting:

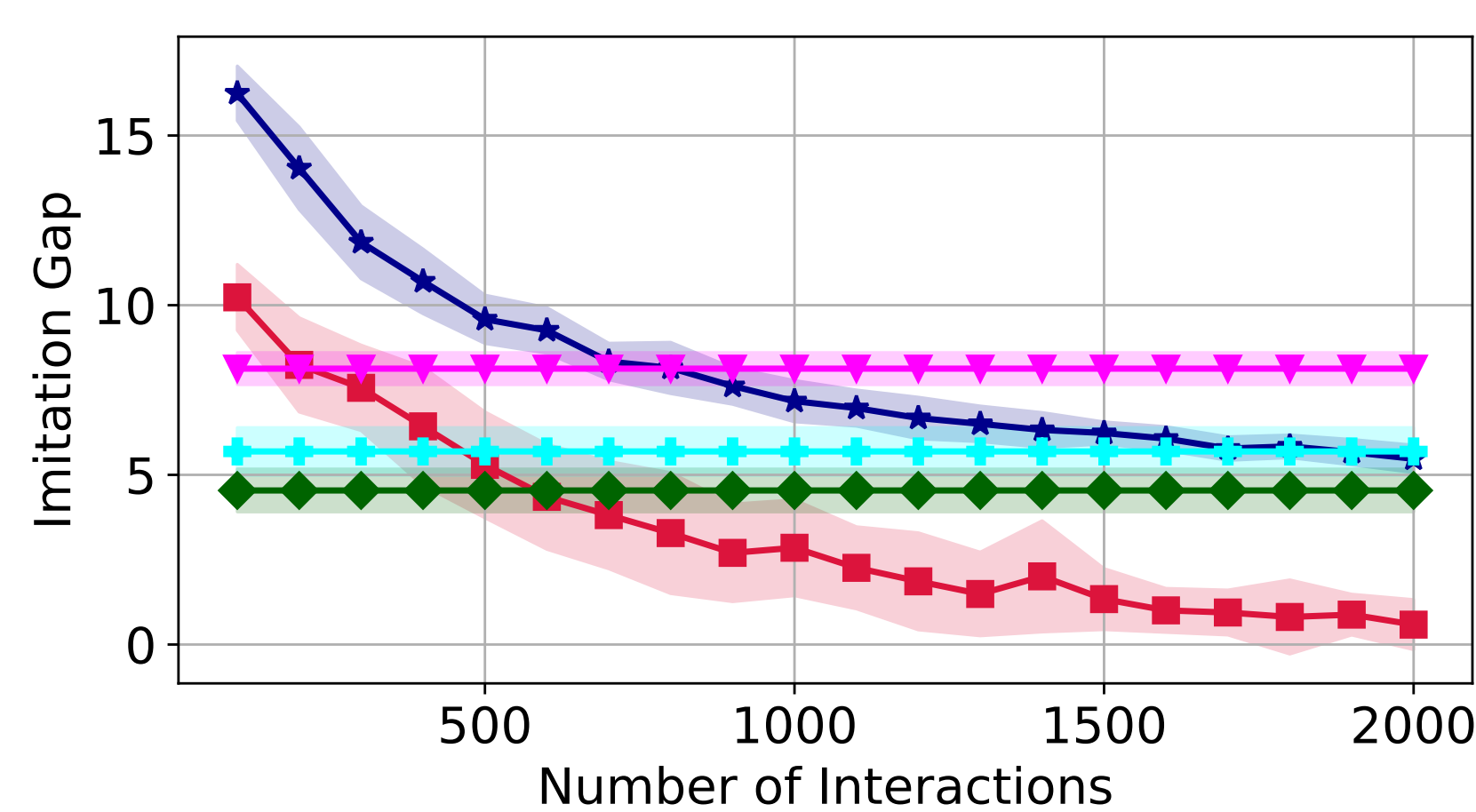
- Unknown π^E : m trajectories (of expert) are given.
- Unknown P : n trajectories (of learner) to interact with env.

Main Results

Theoretical guarantees of IL algorithms with unknown transitions

	Expert Sample Complexity (m)	Interaction Complexity (n)
BC[1]	$\tilde{\mathcal{O}}\left(\frac{H^2 \mathcal{S} }{\varepsilon}\right)$	0
FEM[2]	$\tilde{\mathcal{O}}\left(\frac{H^2 \mathcal{S} }{\varepsilon^2} + \frac{H^8 \mathcal{S} ^3 \mathcal{A} }{\varepsilon^5}\right)$	0
GTAL[3]	$\tilde{\mathcal{O}}\left(\frac{H^2 \mathcal{S} }{\varepsilon^2} + \frac{H^6 \mathcal{S} ^3 \mathcal{A} }{\varepsilon^3}\right)$	0
OAL[4]	$\tilde{\mathcal{O}}\left(\frac{H^2 \mathcal{S} }{\varepsilon^2}\right)$	$\tilde{\mathcal{O}}\left(\frac{H^4 \mathcal{S} ^2 \mathcal{A} }{\varepsilon^2}\right)$
MB-TAIL (this work)	$\tilde{\mathcal{O}}\left(\frac{H^{3/2} \mathcal{S} }{\varepsilon}\right)$	$\tilde{\mathcal{O}}\left(\frac{H^3 \mathcal{S} ^2 \mathcal{A} }{\varepsilon^2}\right)$

— MB-TAIL — OAL — BC — FEM — GTAL



Algorithm Performance

Remark

- Theoretically improvements on both m and n
- Better empirical performance

Algorithmic Framework

Algorithm 1 Meta-algorithm for AIL with Unknown Transitions

Require: Expert demonstrations \mathcal{D} .

- 1: $\hat{P} \leftarrow$ Invoke a method to interact with the env and then learn a transition model.
- 2: $\tilde{d}_h^{\pi^E} \leftarrow$ Estimate the expert state-action distribution from \mathcal{D} .
- 3: $\bar{\pi} \leftarrow$ Apply AIL to perform distribution matching with $\tilde{d}_h^{\pi^E}$ and \hat{P} .

Ensure: Policy $\bar{\pi}$.

Proposition 1: Error Decomposition

(a) \hat{P} is $(\varepsilon_{\text{RFE}}, \delta_{\text{RFE}})$ -PAC for uniform policy evaluation (UPE):

$$\mathbb{P}\left(\text{for any } r, \pi \in \Pi, |\mathbf{V}^{\pi, P, r} - \mathbf{V}^{\pi, \hat{P}, r}| \leq \varepsilon\right) \geq 1 - \delta.$$

(b) $\tilde{d}_h^{\pi^E}$ is $(\varepsilon_{\text{EST}}, \delta_{\text{EST}})$ -PAC for estimating $d_h^{\pi^E}$:

$$\mathbb{P}\left(\sum_{h=1}^H \left\| d_h^{\pi^E} - \tilde{d}_h^{\pi^E} \right\| \leq \varepsilon_{\text{EST}}\right) \geq 1 - \delta_{\text{EST}}.$$

(c) $\bar{\pi}$ is ε_{OPT} -optimal solution:

$$\sum_{h=1}^H \left\| d_h^{\bar{\pi}, \hat{P}} - \widehat{d}_h^{\pi^E} \right\|_1 \leq \min_{\pi \in \Pi} \sum_{h=1}^H \left\| d_h^{\pi, \hat{P}} - \widehat{d}_h^{\pi^E} \right\|_1 + \varepsilon_{\text{OPT}}.$$

Then we have

$$\mathbb{P}(\mathbf{V}^{\pi^E} - \mathbf{V}^{\bar{\pi}} \leq 2\varepsilon_{\text{EST}} + 2\varepsilon_{\text{RFE}} + \varepsilon_{\text{OPT}}) \geq 1 - \delta_{\text{EST}} - \delta_{\text{RFE}}.$$

Proposition 1 characterizes three types of errors in the training of AIL:

$$\text{Imitation Gap} \lesssim \text{Exp. Error} + \text{Est. Error} + \text{Opt. Error}$$

Theoretical Analysis

(Part a) Efficient Exploration for UPE

Definition 1: Reward-free Exploration (RFE)

An algorithm is (ε, δ) -PAC for RFE if

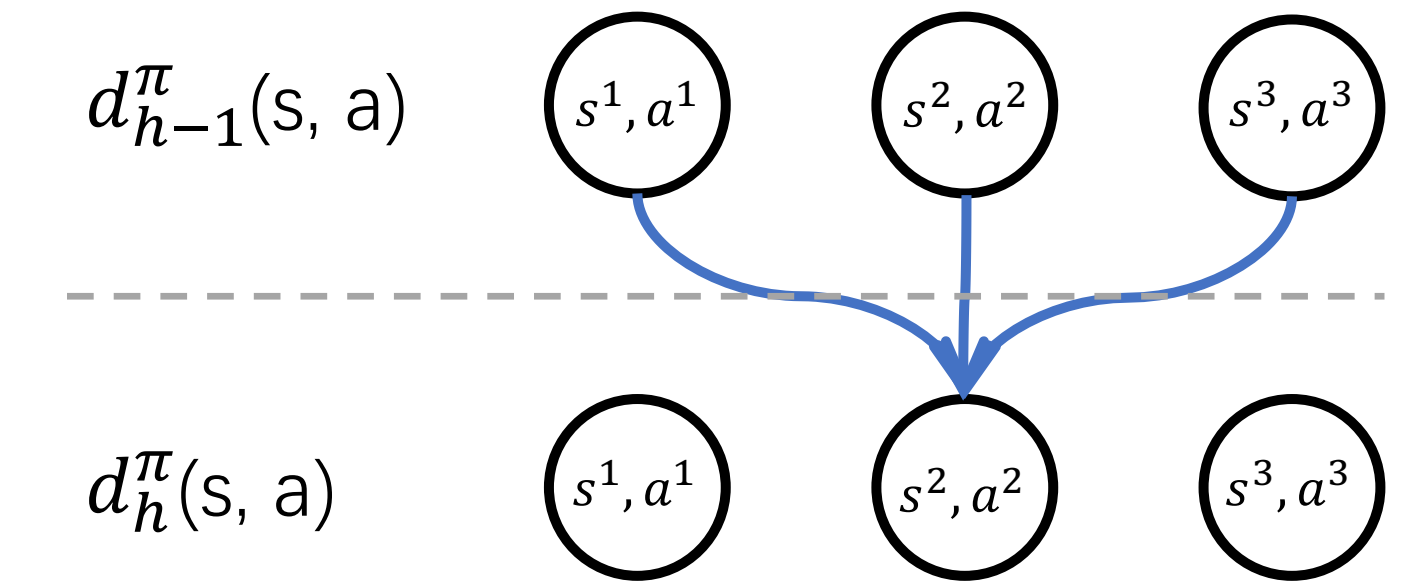
$$\mathbb{P}\left(\text{for any reward function } r, |\mathbf{V}^{\pi^*} - \mathbf{V}^{\hat{\pi}^*}| \leq \varepsilon\right) \geq 1 - \delta.$$

Lemma 1: RF-Express Solves RFE (and UPE)

RF-Express Algorithm in [5] ensures (ε, δ) -PAC for UPE if

$$n \gtrsim \frac{H^3|\mathcal{S}||\mathcal{A}|}{\varepsilon^2} \left(|\mathcal{S}| + \log\left(\frac{|\mathcal{S}|H}{\delta}\right) \right).$$

(Part b) Optimal State-action Distribution Estimation



Bellman-flow equation: $d_h^{\pi^E}(s, a) = (\sum_{s', a'} d_{h-1}^{\pi^E}(s', a') P_{h-1}(s|s', a')) \pi_h^E(a|s)$.

Split \mathcal{D} into two parts: $\mathcal{D} = \mathcal{D}^1 \cup \mathcal{D}_1^c$. Decompose $d_h^{\pi^E}(s, a) =$

$$\sum_{\text{tr}_h \in \text{Tr}_h^{\mathcal{D}^1}} \mathbb{P}^{\pi^E}(\text{tr}_h) \mathbb{I}\{\text{tr}_h(\cdot, \cdot) = (s, a)\} + \sum_{\text{tr}_h \notin \text{Tr}_h^{\mathcal{D}^1}} \mathbb{P}^{\pi^E}(\text{tr}_h) \mathbb{I}\{\text{tr}_h(\cdot, \cdot) = (s, a)\}$$

Here $\text{Tr}_h^{\mathcal{D}^1}$ is the set of sub-trajectories along which each state is covered in \mathcal{D}^1 . Estimate \spadesuit from the complementary set \mathcal{D}_1^c :

$$\frac{\sum_{\text{tr}_h \in \mathcal{D}_1^c} \mathbb{I}\{\text{tr}_h(\cdot, \cdot) = (s, a), \text{tr}_h \notin \text{Tr}_h^{\mathcal{D}^1}\}}{|\mathcal{D}_1^c|}$$

Estimate \clubsuit from trajectories $\mathcal{D}'_{\text{env}}$ collected by the BC's policy:

$$\frac{\sum_{\text{tr}_h \in \mathcal{D}'_{\text{env}}} \mathbb{I}\{\text{tr}_h(\cdot, \cdot) = (s, a), \text{tr}_h \in \text{Tr}_h^{\mathcal{D}^1}\}}{|\mathcal{D}'_{\text{env}}|}$$

Lemma 2

The estimator $\tilde{d}_h^{\pi^E}$ is (ε, δ) -PAC for estimating $d_h^{\pi^E}$ if

$$m \gtrsim \frac{H^{3/2}|\mathcal{S}|}{\varepsilon} \log\left(\frac{|\mathcal{S}|H}{\delta}\right), \quad n' \gtrsim \frac{H^2|\mathcal{S}|}{\varepsilon^2} \log\left(\frac{|\mathcal{S}|H}{\delta}\right).$$

(Part c) Efficient Optimization

Minimax formulation of state-action distribution matching

$$\max_{w \in \mathcal{W}} \min_{\pi \in \Pi} \sum_{h=1}^H \sum_{(s, a)} w_h(s, a) (\tilde{d}_h^{\pi^E}(s, a) - d_h^{\pi, \hat{P}}(s, a)).$$

Apply online optimization methods approximately to obtain the saddle point.

- $\pi^{(t)} \leftarrow$ Solve the RL optimal policy with \hat{P} and $w^{(t)}$.
- Update $w^{(t+1)} := \mathcal{P}_{\mathcal{W}}(w^{(t)} - \eta^{(t)} \nabla f^{(t)}(w^{(t)}))$.

Lemma 3

The optimization procedure can return an ε -accurate solution if $T \gtrsim \frac{H^2|\mathcal{S}||\mathcal{A}|}{\varepsilon^2}$.

References:

- [1] Rajaraman, N., et al. "Toward the fundamental limits of imitation learning." NeurIPS, 2020.
- [2] Abbeel, P., et al. "Apprenticeship learning via inverse reinforcement learning." ICML, 2004.
- [3] Syed, U., et al. "A game-theoretic approach to apprenticeship learning." NeurIPS, 2007.
- [4] Shani, L., et al. "Online apprenticeship learning." AAAI, 2022.
- [5] Menard, P., et al. "Fast active learning for pure exploration in reinforcement learning." ICML, 2021.

Code at GitHub:



Paper at arXiv:

