



POLITECNICO
MILANO 1863



Towards Theoretical Understanding of Inverse Reinforcement Learning

Alberto Maria Metelli¹, Filippo Lazzati¹ and Marcello Restelli¹

¹ Politecnico di Milano

40th International Conference on Machine Learning, Honolulu, HI.

July 2023

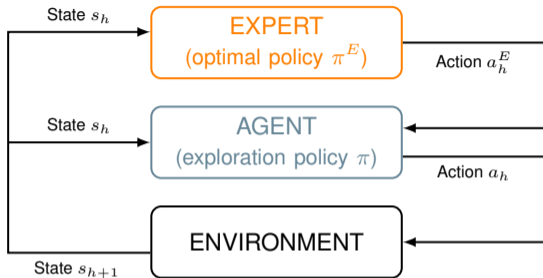
Research Question

How many samples are needed to accurately solve the Inverse Reinforcement Learning (IRL) problem with high probability?

Research Question

How many samples are needed to accurately solve the Inverse Reinforcement Learning (IRL) problem with high probability?

Sample Complexity **Lower Bound** for IRL



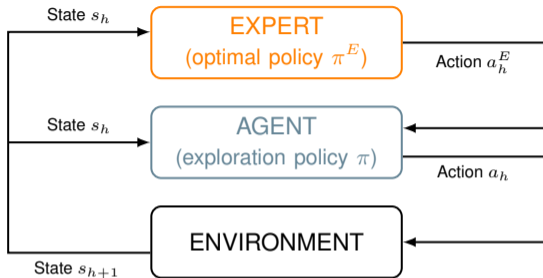
■ At every stage $h \in \llbracket H \rrbracket$:

- Observe state s_h
- Observe **expert action**
 $a_h^E \sim \pi_h^E(\cdot | s_h)$
- Play **exploratory action**
 $a_h \sim \pi_h(\cdot | s_h)$
- Transition to next state
 $s_{h+1} \sim p_h(\cdot | s_h, a_h)$

■ Traditional Goal of IRL (Arora and Doshi, 2021; Adams et al., 2022)

Find one feasible reward function r^* that makes the expert's policy π^E optimal, i.e.,

$$\pi^E \in \arg \max_{\pi} V^{\pi}(\cdot; r^*)$$



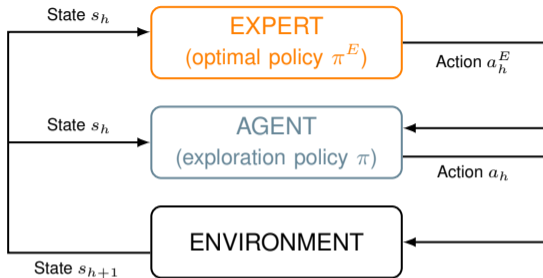
■ At every stage $h \in \llbracket H \rrbracket$:

- Observe state s_h
- Observe **expert action**
 $a_h^E \sim \pi_h^E(\cdot | s_h)$
- Play **exploratory action**
 $a_h \sim \pi_h(\cdot | s_h)$
- Transition to next state
 $s_{h+1} \sim p_h(\cdot | s_h, a_h)$

■ Traditional Goal of IRL (Arora and Doshi, 2021; Adams et al., 2022)

Find one feasible reward function r^* that makes the expert's policy π^E optimal, i.e.,

$$\pi^E \in \arg \max_{\pi} V^{\pi}(\cdot; r^*)$$

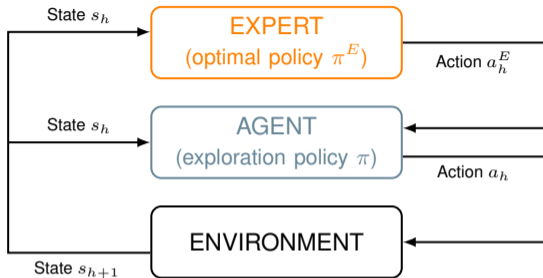


- At every stage $h \in \llbracket H \rrbracket$:
 - Observe state s_h
 - Observe **expert action**
 $a_h^E \sim \pi_h^E(\cdot | s_h)$
 - Play **exploratory action**
 $a_h \sim \pi_h(\cdot | s_h)$
 - Transition to next state
 $s_{h+1} \sim p_h(\cdot | s_h, a_h)$

- Traditional Goal of IRL (Arora and Doshi, 2021; Adams et al., 2022)

Find one feasible reward function r^* that makes the expert's policy π^E optimal, i.e.,

$$\pi^E \in \arg \max_{\pi} V^{\pi}(\cdot; r^*)$$

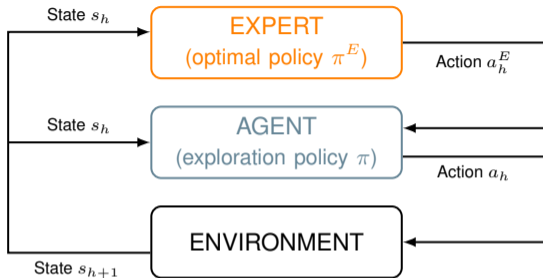


- At every stage $h \in \llbracket H \rrbracket$:
 - Observe state s_h
 - Observe **expert action**
 $a_h^E \sim \pi_h^E(\cdot | s_h)$
 - Play **exploratory action**
 $a_h \sim \pi_h(\cdot | s_h)$
 - Transition to next state
 $s_{h+1} \sim p_h(\cdot | s_h, a_h)$

- Traditional Goal of IRL (Arora and Doshi, 2021; Adams et al., 2022)

Find one feasible reward function r^* that makes the expert's policy π^E optimal, i.e.,

$$\pi^E \in \arg \max_{\pi} V^{\pi}(\cdot; r^*)$$

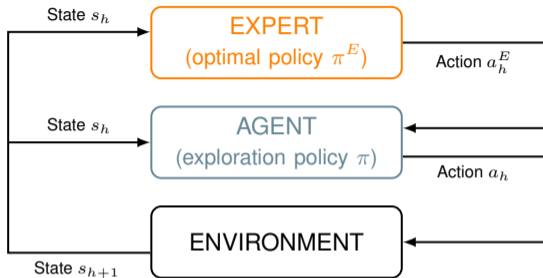


- At every stage $h \in \llbracket H \rrbracket$:
 - Observe state s_h
 - Observe **expert action** $a_h^E \sim \pi_h^E(\cdot | s_h)$
 - Play **exploratory action** $a_h \sim \pi_h(\cdot | s_h)$
 - Transition to next state $s_{h+1} \sim p_h(\cdot | s_h, a_h)$

- Traditional Goal of IRL (Arora and Doshi, 2021; Adams et al., 2022)

Find **one feasible reward function** r^* that makes the **expert's policy** π^E **optimal**, i.e.,

$$\pi^E \in \arg \max_{\pi} V^{\pi}(\cdot; r^*)$$



- At every stage $h \in \llbracket H \rrbracket$:
 - Observe state s_h
 - Observe **expert action** $a_h^E \sim \pi_h^E(\cdot | s_h)$
 - Play **exploratory action** $a_h \sim \pi_h(\cdot | s_h)$
 - Transition to next state $s_{h+1} \sim p_h(\cdot | s_h, a_h)$

- Traditional Goal of IRL (Arora and Doshi, 2021; Adams et al., 2022)

Find **one feasible reward function** r^* that makes the **expert's policy** π^E **optimal**, i.e.,

$$\pi^E \in \arg \max_{\pi} V^{\pi}(\cdot; r^*)$$

- **Ambiguity** problem (Ng and Russell, 2000)
- Study the **full** set of feasible rewards r^* \rightarrow **Feasible Reward Set** (Metelli et al., 2021)

Find all feasible reward functions \mathcal{R} that make the expert's policy π^E optimal, i.e.,

$$\mathcal{R} := \left\{ \text{all rewards } r^* : \pi^E \in \arg \max_{\pi} V^{\pi}(\cdot; r^*) \right\}$$

- \mathcal{R} defined through **linear constraints**
- If p and π^E are **known**, check if \hat{r} is feasible takes $O(HS^2A)$

- **Ambiguity** problem (Ng and Russell, 2000)
- Study the **full** set of feasible rewards r^* → **Feasible Reward Set** (Metelli et al., 2021)

Find **all feasible reward functions** \mathcal{R} that make the **expert's policy** π^E **optimal**, i.e.,

$$\mathcal{R} := \left\{ \text{all rewards } r^* : \pi^E \in \arg \max_{\pi} V^{\pi}(\cdot; r^*) \right\}$$

- \mathcal{R} defined through **linear constraints**
- If p and π^E are **known**, check if \hat{r} is feasible takes $O(HS^2A)$

- **Ambiguity** problem (Ng and Russell, 2000)
- Study the **full** set of feasible rewards r^* \rightarrow **Feasible Reward Set** (Metelli et al., 2021)

Find **all feasible reward functions** \mathcal{R} that make the **expert's policy** π^E **optimal**, i.e.,

$$\mathcal{R} := \left\{ \text{all rewards } r^* : \pi^E \in \arg \max_{\pi} V^{\pi}(\cdot; r^*) \right\}$$

- \mathcal{R} defined through **linear constraints**
- If p and π^E are **known**, check if \hat{r} is feasible takes $O(HS^2A)$

- **Ambiguity** problem (Ng and Russell, 2000)
- Study the **full** set of feasible rewards r^* \rightarrow **Feasible Reward Set** (Metelli et al., 2021)

Find **all feasible reward functions** \mathcal{R} that make the **expert's policy** π^E **optimal**, i.e.,

$$\mathcal{R} := \left\{ \text{all rewards } r^* : \pi^E \in \arg \max_{\pi} V^{\pi}(\cdot; r^*) \right\}$$

- \mathcal{R} defined through **linear constraints**
- If p and π^E are **known**, check if \hat{r} is feasible takes $O(HS^2A)$

- **Ambiguity** problem (Ng and Russell, 2000)
- Study the **full** set of feasible rewards r^* \rightarrow **Feasible Reward Set** (Metelli et al., 2021)

Find **all feasible reward functions** \mathcal{R} that make the **expert's policy** π^E **optimal**, i.e.,

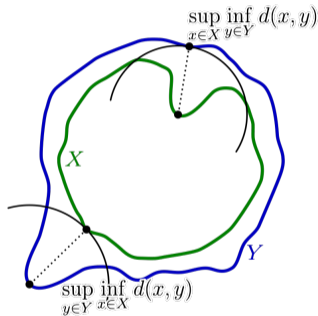
$$\mathcal{R} := \left\{ \text{all rewards } r^* : \pi^E \in \arg \max_{\pi} V^{\pi}(\cdot; r^*) \right\}$$

- \mathcal{R} defined through **linear constraints**
- If p and π^E are **known**, check if \hat{r} is feasible takes $O(HS^2A)$

- Transition model p and expert's policy π^E **unknown**
- Estimate \hat{p} and $\hat{\pi}^E$ with **samples** inducing $\hat{\mathcal{R}}$
- Hausdorff distance** between \mathcal{R} (true) and $\hat{\mathcal{R}}$ (estimated) feasible reward sets

$$\mathcal{H}_d(\mathcal{R}, \hat{\mathcal{R}}) = \max \left\{ \sup_{r \in \mathcal{R}} \inf_{\hat{r} \in \hat{\mathcal{R}}} d(r, \hat{r}), \sup_{\hat{r} \in \hat{\mathcal{R}}} \inf_{r \in \mathcal{R}} d(r, \hat{r}) \right\}$$

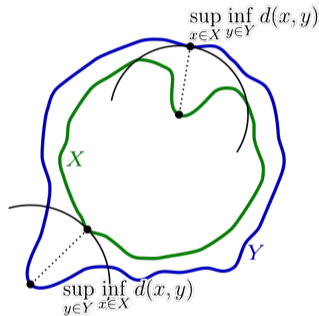
where
$$d(r, \hat{r}) = \max_{s, a, h} |r_h(s, a) - \hat{r}_h(s, a)|$$



- Transition model p and expert's policy π^E **unknown**
- Estimate \hat{p} and $\hat{\pi}^E$ with **samples** inducing $\hat{\mathcal{R}}$
- Hausdorff distance** between \mathcal{R} (true) and $\hat{\mathcal{R}}$ (estimated) feasible reward sets

$$\mathcal{H}_d(\mathcal{R}, \hat{\mathcal{R}}) = \max \left\{ \sup_{r \in \mathcal{R}} \inf_{\hat{r} \in \hat{\mathcal{R}}} d(r, \hat{r}), \sup_{\hat{r} \in \hat{\mathcal{R}}} \inf_{r \in \mathcal{R}} d(r, \hat{r}) \right\}$$

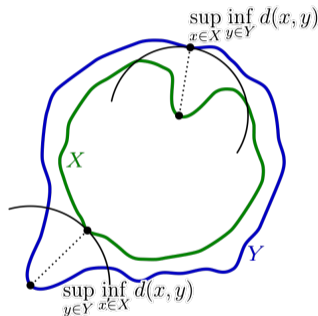
where
$$d(r, \hat{r}) = \max_{s, a, h} |r_h(s, a) - \hat{r}_h(s, a)|$$



- Transition model p and expert's policy π^E **unknown**
- Estimate \hat{p} and $\hat{\pi}^E$ with **samples** inducing $\hat{\mathcal{R}}$
- Hausdorff distance** between \mathcal{R} (true) and $\hat{\mathcal{R}}$ (estimated) feasible reward sets

$$\mathcal{H}_d(\mathcal{R}, \hat{\mathcal{R}}) = \max \left\{ \sup_{r \in \mathcal{R}} \inf_{\hat{r} \in \hat{\mathcal{R}}} d(r, \hat{r}), \sup_{\hat{r} \in \hat{\mathcal{R}}} \inf_{r \in \mathcal{R}} d(r, \hat{r}) \right\}$$

where
$$d(r, \hat{r}) = \max_{s, a, h} |r_h(s, a) - \hat{r}_h(s, a)|$$



- An algorithm \mathfrak{A} that outputs $\hat{\mathcal{R}}_\tau$ after τ calls to the environment is (ϵ, δ) -PAC if

$$\mathbb{P}\left(\mathcal{H}_d\left(\mathcal{R}, \hat{\mathcal{R}}_\tau\right) > \epsilon\right) \leq \delta$$

$\tau =$ sample complexity

- Sample complexity lower bound

$$\tau \geq \text{poly}\left(S, A, H, \frac{1}{\epsilon}, \log\left(\frac{1}{\delta}\right)\right)$$

- An algorithm \mathfrak{A} that outputs $\hat{\mathcal{R}}_\tau$ after τ calls to the environment is (ϵ, δ) -PAC if

$$\mathbb{P}\left(\mathcal{H}_d\left(\mathcal{R}, \hat{\mathcal{R}}_\tau\right) > \epsilon\right) \leq \delta$$

$\tau =$ sample complexity

- **Sample complexity lower bound**

$$\tau \geq \text{poly}\left(S, A, H, \frac{1}{\epsilon}, \log\left(\frac{1}{\delta}\right)\right)$$

Theorem

For any (ϵ, δ) -PAC algorithm \mathfrak{A} , with ϵ and δ sufficiently small, there exists an IRL problem, with S , A and H sufficiently large, such that the **expected sample complexity** is lower bounded by:

- if the transition model p is time-inhomogeneous (i.e., $p_h \neq p_{h+1}$):

$$\mathbb{E}[\tau] \geq \Omega\left(\frac{H^3 SA}{\epsilon^2} \left(\log\left(\frac{1}{\delta}\right) + S\right)\right);$$

- if the transition model p is time-homogeneous (i.e., $p_h = p_{h+1}$):

$$\mathbb{E}[\tau] \geq \Omega\left(\frac{H^2 SA}{\epsilon^2} \left(\log\left(\frac{1}{\delta}\right) + S\right)\right).$$

Theorem

For any (ϵ, δ) -PAC algorithm \mathfrak{A} , with ϵ and δ sufficiently small, there exists an IRL problem, with S , A and H sufficiently large, such that the **expected sample complexity** is lower bounded by:

- if the transition model p is time-inhomogeneous (i.e., $p_h \neq p_{h+1}$):

$$\mathbb{E}[\tau] \geq \Omega \left(\frac{H^3 SA}{\epsilon^2} \left(\log \left(\frac{1}{\delta} \right) + S \right) \right);$$

- if the transition model p is time-homogeneous (i.e., $p_h = p_{h+1}$):

$$\mathbb{E}[\tau] \geq \Omega \left(\frac{H^2 SA}{\epsilon^2} \left(\log \left(\frac{1}{\delta} \right) + S \right) \right).$$

Theorem

For any (ϵ, δ) -PAC algorithm \mathfrak{A} , with ϵ and δ sufficiently small, there exists an IRL problem, with S , A and H sufficiently large, such that the **expected sample complexity** is lower bounded by:

- if the transition model p is time-inhomogeneous (i.e., $p_h \neq p_{h+1}$):

$$\mathbb{E}[\tau] \geq \Omega\left(\frac{H^3 SA}{\epsilon^2} \left(\log\left(\frac{1}{\delta}\right) + S\right)\right);$$

- if the transition model p is time-homogeneous (i.e., $p_h = p_{h+1}$):

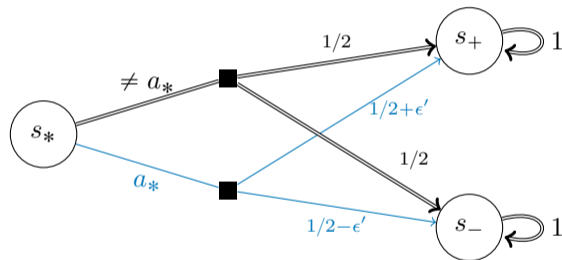
$$\mathbb{E}[\tau] \geq \Omega\left(\frac{H^2 SA}{\epsilon^2} \left(\log\left(\frac{1}{\delta}\right) + S\right)\right).$$

Two regimes of δ

Two **regimes** of $\delta \rightarrow$ **Small- δ regime**

- Expert's policy $\pi^E(s) = a_0$
- Hard to identify which action behaves like a_*
- Construct $\Theta(A)$ hard instances
- Technical tool:** Bretagnolle-Huber inequality (Lattimore and Szepesvári, 2020)

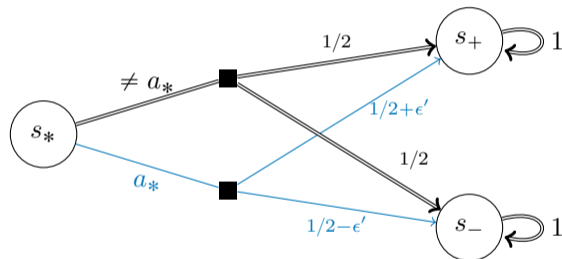
$$\Omega\left(\frac{H^3 SA}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right)$$



Two **regimes** of $\delta \rightarrow$ **Small- δ regime**

- Expert's policy $\pi^E(s) = a_0$
- **Hard** to identify which action behaves like a_*
- Construct $\Theta(A)$ hard instances
- **Technical tool**: Bretagnolle-Huber inequality (Lattimore and Szepesvári, 2020)

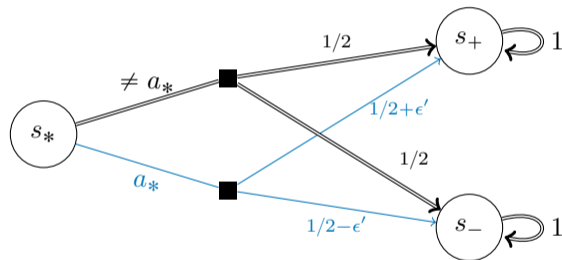
$$\Omega\left(\frac{H^3 SA}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right)$$



Two **regimes** of $\delta \rightarrow$ **Small- δ regime**

- Expert's policy $\pi^E(s) = a_0$
- **Hard** to identify which action behaves like a_*
- Construct $\Theta(A)$ hard instances
- **Technical tool**: Bretagnolle-Huber inequality (Lattimore and Szepesvári, 2020)

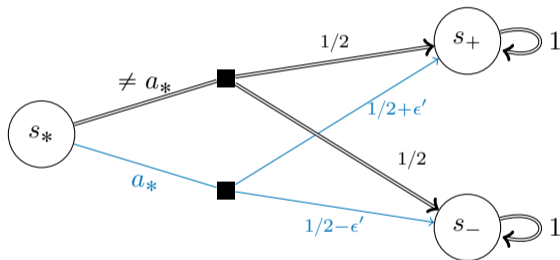
$$\Omega\left(\frac{H^3 SA}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right)$$



Two **regimes** of $\delta \rightarrow$ **Small- δ regime**

- Expert's policy $\pi^E(s) = a_0$
- Hard** to identify which action behaves like a_*
- Construct $\Theta(A)$ hard instances
- Technical tool**: Bretagnolle-Huber inequality (Lattimore and Szepesvári, 2020)

$$\Omega\left(\frac{H^3 SA}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right)$$

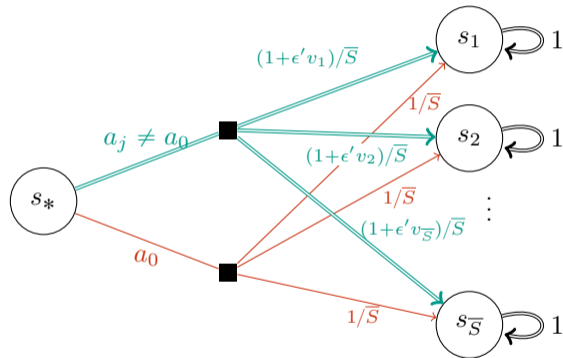


Two regimes of δ

Two regimes of $\delta \rightarrow$ **Large- δ regime**

- Expert's policy $\pi^E(s) = a_0$
- Hard to distinguish all actions a_j
- Construct $\Theta(2^S)$ hard instances via a packing argument based on Hamming coding
- Technical tool: Fano's inequality (Gerchinovitz et al., 2020)

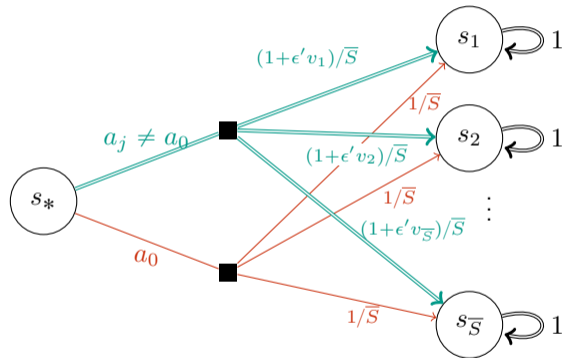
$$\Omega\left(\frac{H^3 S^2 A}{\epsilon^2}\right)$$



Two regimes of $\delta \rightarrow$ **Large- δ regime**

- Expert's policy $\pi^E(s) = a_0$
- Hard** to distinguish all actions a_j
- Construct $\Theta(2^S)$ hard instances via a **packing argument** based on **Hamming coding**
- Technical tool**: Fano's inequality (Gerchinovitz et al., 2020)

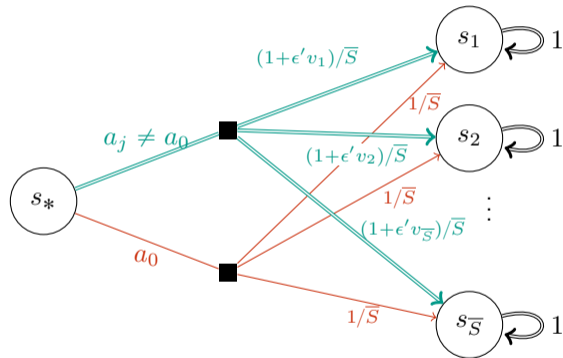
$$\Omega\left(\frac{H^3 S^2 A}{\epsilon^2}\right)$$



Two regimes of $\delta \rightarrow$ **Large- δ regime**

- Expert's policy $\pi^E(s) = a_0$
- Hard** to distinguish all actions a_j
- Construct $\Theta(2^S)$ hard instances via a **packing** argument based on **Hamming coding**
- Technical tool: Fano's inequality (Gerchinovitz et al., 2020)

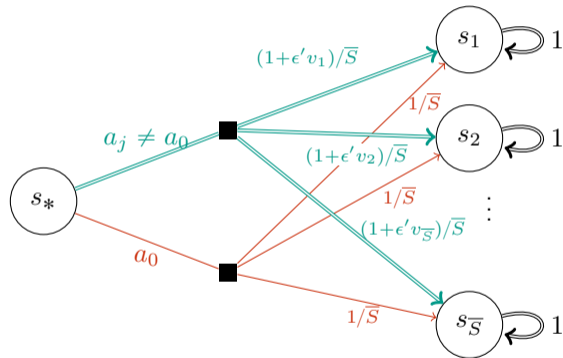
$$\Omega\left(\frac{H^3 S^2 A}{\epsilon^2}\right)$$



Two regimes of $\delta \rightarrow$ **Large- δ regime**

- Expert's policy $\pi^E(s) = a_0$
- Hard** to distinguish all actions a_j
- Construct $\Theta(2^S)$ hard instances via a **packing** argument based on **Hamming coding**
- Technical tool**: Fano's inequality (Gerchinovitz et al., 2020)

$$\Omega\left(\frac{H^3 S^2 A}{\epsilon^2}\right)$$



- **Lipschitz** properties of the feasible reward set \mathcal{R}
- **Uniform sampling** algorithm nearly matches the lower bound
- Relation with different **dissimilarity index** between rewards

- **Lipschitz** properties of the feasible reward set \mathcal{R}
- **Uniform sampling** algorithm nearly matches the lower bound
- Relation with different **dissimilarity index** between rewards

- **Lipschitz** properties of the feasible reward set \mathcal{R}
- **Uniform sampling** algorithm nearly matches the lower bound
- Relation with different **dissimilarity index** between rewards

**Thank You
for Your
Attention!**



Contacts: albertomaria.metelli@polimi.it
Link: <https://icml.cc/virtual/2023/poster/24193>

- Adams, S. C., Cody, T., and Beling, P. A. (2022). A survey of inverse reinforcement learning. *Artif. Intell. Rev.*, 55(6):4307–4346.
- Arora, S. and Doshi, P. (2021). A survey of inverse reinforcement learning: Challenges, methods and progress. *Artif. Intell.*, 297:103500.
- Dann, C. and Brunskill, E. (2015). Sample complexity of episodic fixed-horizon reinforcement learning. In *Advances in Neural Information Processing Systems 28 (NIPS)*, pages 2818–2826.
- Domingues, O. D., Ménard, P., Kaufmann, E., and Valko, M. (2021). Episodic reinforcement learning in finite mdps: Minimax lower bounds revisited. In *Algorithmic Learning Theory (ALT)*, volume 132 of *Proceedings of Machine Learning Research*, pages 578–598. PMLR.
- Gerchinovitz, S., Ménard, P., and Stoltz, G. (2020). Fano’s inequality for random variables. *Statistical Science*, 35(2):178–201.
- Jin, C., Krishnamurthy, A., Simchowitz, M., and Yu, T. (2020). Reward-free exploration for reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, volume 119 of *Proceedings of Machine Learning Research*, pages 4870–4879. PMLR.
- Lattimore, T. and Szepesvári, C. (2020). *Bandit algorithms*. Cambridge University Press.
- Metelli, A. M., Ramponi, G., Concetti, A., and Restelli, M. (2021). Provably efficient learning of transferable rewards. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, volume 139 of *Proceedings of Machine Learning Research*, pages 7665–7676. PMLR.
- Ng, A. Y. and Russell, S. (2000). Algorithms for inverse reinforcement learning. In *Proceedings of the Seventeenth International Conference on Machine Learning (ICML)*, pages 663–670. Morgan Kaufmann.