

Active coverage for PAC Reinforcement Learning

Aymen Al Marjani^{1,2}, Andrea Tirinzoni³, Emilie Kaufmann^{2,3},

¹UMPA, ENS de Lyon

²Inria Lille, Equipe SCOOOL

³CNRS, CRISAL ⁴META AI, Paris

28 Avril 2023



Outline

- 1 Framework & definition of active coverage
- 2 Why is it worth studying ?
- 3 Sample complexity lower bound
- 4 Vanilla CovGame: a simple sometimes-optimal algorithm
- 5 Improved CovGame: a (near-)optimal algorithm
- 6 Discussion

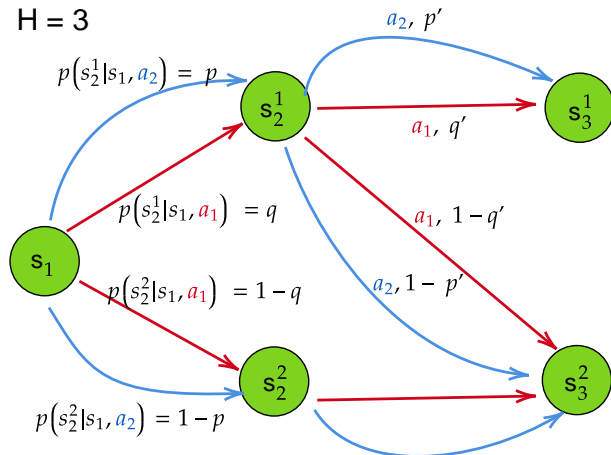
Outline

- 1 Framework & definition of active coverage
- 2 Why is it worth studying ?
- 3 Sample complexity lower bound
- 4 Vanilla CovGame: a simple sometimes-optimal algorithm
- 5 Improved CovGame: a (near-)optimal algorithm
- 6 Discussion

Finite-horizon Tabular MDPs

$$\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, \{p_h\}_{h \in [H]}, s_1)$$

$H = 3$



$$p(s_3^2 | s_2^2, a_1) = p(s_3^2 | s_2^2, a_2) = 1$$

Active coverage: learning problem

The Active Coverage problem: Design an algorithm \mathbb{A} which takes as input a target function $c : [H] \times \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}_+$ and a failure probability δ . \mathbb{A} should interact with \mathcal{M} and return a dataset \mathcal{D}_t such that

$$\mathbb{P}_{\mathcal{M}, \mathbb{A}} \left(\exists t \geq 1, \forall (h, s, a), n_h(s, a; \mathcal{D}_t) \geq c_h(s, a) \right) \geq 1 - \delta.$$

Active coverage: learning problem

The Active Coverage problem: Design an algorithm \mathbb{A} which takes as input a target function $c : [H] \times \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}_+$ and a failure probability δ . \mathbb{A} should interact with \mathcal{M} and return a dataset \mathcal{D}_t such that

$$\mathbb{P}_{\mathcal{M}, \mathbb{A}} \left(\exists t \geq 1, \forall (h, s, a), n_h(s, a; \mathcal{D}_t) \geq c_h(s, a) \right) \geq 1 - \delta.$$

Performance criteria:

$$\tau := \inf \{ t \geq 1, \forall (h, s, a), n_h(s, a; \mathcal{D}_t) \geq c_h(s, a) \}$$

\implies We want \mathbb{A} to make τ as small as possible.

Active coverage: Protocol of interaction

Algorithm 1 Protocol of interaction

- 1: **Input:** target function c , risk $\delta \in (0, 1)$.
 - 2: Initialize dataset $\mathcal{D}_0 \leftarrow \emptyset$
 - 3: Set target set $\mathcal{X} = \{(s, a, h) \in [H] \times \mathcal{S} \times \mathcal{A} : c_h(s, a) > 0\}$
 - 4: **for** $t = 1, 2, \dots$ **do**
 - 5: $\pi^t \leftarrow \text{COVERAGEALGORITHM}()$
 - 6: Play π^t and observe trajectory $\mathcal{H}_t := \{(s_h^t, a_h^t, s_{h+1}^t)\}_{1 \leq h \leq H-1}$
 - 7: Update dataset $\mathcal{D}_t \leftarrow \mathcal{D}_{t-1} \cup \mathcal{H}_t$.
 - 8: **If** $\forall (h, s, a), n_h(s, a; \mathcal{D}_t) \geq c_h(s, a)$:
 - 9: Stop and return \mathcal{D}_t
 - 10: **end for**
-

where the counts $n_h(s, a; \mathcal{D}_t) = \sum_{u=1}^t \mathbb{1}(s_h^u = s, a_h^u = a)$.

Outline

- 1 Framework & definition of active coverage
- 2 Why is it worth studying ?**
- 3 Sample complexity lower bound
- 4 Vanilla CovGame: a simple sometimes-optimal algorithm
- 5 Improved CovGame: a (near-)optimal algorithm
- 6 Discussion

Applications

- Best-Policy Identification: Given some reward function $r : [H] \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, find an ε -optimal policy, i.e. $V_1^{\hat{\pi}}(s_1) \geq V_1^{\pi^*}(s_1) - \varepsilon$.
- Reward-free Exploration: Collect Dataset \mathcal{D} such that you can find ε -optimal policy for any reward function r using \mathcal{D} .
- Hybrid RL (Offline + Online): Complete some dataset \mathcal{D} that was given to you.

Outline

- 1 Framework & definition of active coverage
- 2 Why is it worth studying ?
- 3 Sample complexity lower bound**
- 4 Vanilla CovGame: a simple sometimes-optimal algorithm
- 5 Improved CovGame: a (near-)optimal algorithm
- 6 Discussion

Lower bound: Intuition

Imagine that the agent chooses to explore using the policy π^{exp} , for all episodes.

Let $p_h^{\pi^{exp}}(s, a) = \mathbb{P}(s_h = s, a_h = a | s_1, \pi^{exp})$ be the probability of visiting (h, s, a) in one episode.

Lower bound: Intuition

Imagine that the agent chooses to explore using the policy π^{exp} , for all episodes.

Let $p_h^{\pi^{exp}}(s, a) = \mathbb{P}(s_h = s, a_h = a | s_1, \pi^{exp})$ be the probability of visiting (h, s, a) in one episode.

The expected number of episodes before the first visit is $1/p_h^{\pi^{exp}}(s, a)$.

Lower bound: Intuition

Imagine that the agent chooses to explore using the policy π^{exp} , for all episodes.

Let $p_h^{\pi^{exp}}(s, a) = \mathbb{P}(s_h = s, a_h = a | s_1, \pi^{exp})$ be the probability of visiting (h, s, a) in one episode.

The expected number of episodes before the first visit is $1/p_h^{\pi^{exp}}(s, a)$.

To satisfy $n_h^\tau(s, a) \geq c_h(s, a)$, the agent would need at least $\max_{h,s,a} \frac{c_h(s,a)}{p_h^{\pi^{exp}}(s,a)}$

Lower bound: Intuition

Imagine that the agent chooses to explore using the policy π^{exp} , for all episodes.

Let $p_h^{\pi^{\text{exp}}}(s, a) = \mathbb{P}(s_h = s, a_h = a | s_1, \pi^{\text{exp}})$ be the probability of visiting (h, s, a) in one episode.

The expected number of episodes before the first visit is $1/p_h^{\pi^{\text{exp}}}(s, a)$.

To satisfy $n_h^\tau(s, a) \geq c_h(s, a)$, the agent would need at least $\max_{h,s,a} \frac{c_h(s, a)}{p_h^{\pi^{\text{exp}}}(s, a)}$

Since the agent has the choice of which policy to play,

$$\mathbb{E}[\tau] \gtrsim \inf_{\pi^{\text{exp}} \in \Pi^S} \max_{h,s,a} \frac{c_h(s, a)}{p_h^{\pi^{\text{exp}}}(s, a)}.$$

Lower bound: Statement

Theorem 1

For any target function c and $\delta \in (0, 1)$, the stopping time τ of any δ -correct c -coverage algorithm satisfies $\mathbb{E}[\tau] \geq (1 - \delta)\varphi^*(c)$, where

$$\varphi^*(c) = \inf_{\pi^{\text{exp}} \in \Pi^{\text{S}}} \max_{(s, a, h) \in \mathcal{X}} \frac{c_h(s, a)}{p_h^{\pi^{\text{exp}}}(s, a)},$$

with $\mathcal{X} := \{(h, s, a) : c_h(s, a) > 0\}$.

Outline

- 1 Framework & definition of active coverage
- 2 Why is it worth studying ?
- 3 Sample complexity lower bound
- 4 Vanilla CovGame: a simple sometimes-optimal algorithm**
- 5 Improved CovGame: a (near-)optimal algorithm
- 6 Discussion

Vanilla CovGame: Intuition

$$\frac{1}{\varphi^*(c)} = \sup_{\pi^{\text{exp}} \in \Pi^S} \min_{(s,a,h) \in \mathcal{X}} \frac{p_h^{\pi^{\text{exp}}}(s,a)}{c_h(s,a)}$$

Vanilla CovGame: Intuition

$$\begin{aligned}\frac{1}{\varphi^*(c)} &= \sup_{\pi^{\text{exp}} \in \Pi^S} \min_{(s,a,h) \in \mathcal{X}} \frac{p_h^{\pi^{\text{exp}}}(s,a)}{c_h(s,a)} \\ &= \sup_{\pi^{\text{exp}} \in \Pi^S} \inf_{\lambda \in \Delta_{\mathcal{X}}} \sum_{h,s,a} \frac{p_h^{\pi^{\text{exp}}}(s,a) \lambda_h(s,a)}{c_h(s,a)}\end{aligned}$$

with $\mathcal{X} := \{(h,s,a) : c_h(s,a) > 0\}$ and $\Delta_{\mathcal{X}}$ = the simplex over \mathcal{X} .

Vanilla CovGame: Intuition

$$\begin{aligned}\frac{1}{\varphi^*(c)} &= \sup_{\pi^{\text{exp}} \in \Pi^S} \min_{(s,a,h) \in \mathcal{X}} \frac{p_h^{\pi^{\text{exp}}}(s,a)}{c_h(s,a)} \\ &= \sup_{\pi^{\text{exp}} \in \Pi^S} \inf_{\lambda \in \Delta_{\mathcal{X}}} \sum_{h,s,a} \frac{p_h^{\pi^{\text{exp}}}(s,a) \lambda_h(s,a)}{c_h(s,a)} \\ &= \text{A matrix game !}\end{aligned}$$

with $\mathcal{X} := \{(h, s, a) : c_h(s, a) > 0\}$ and $\Delta_{\mathcal{X}}$ = the simplex over \mathcal{X} .

Vanilla CovGame: Intuition

$$\begin{aligned}\frac{1}{\varphi^*(c)} &= \sup_{\pi^{\text{exp}} \in \Pi^S} \min_{(s,a,h) \in \mathcal{X}} \frac{p_h^{\pi^{\text{exp}}}(s,a)}{c_h(s,a)} \\ &= \sup_{\pi^{\text{exp}} \in \Pi^S} \inf_{\lambda \in \Delta_{\mathcal{X}}} \sum_{h,s,a} \frac{p_h^{\pi^{\text{exp}}}(s,a) \lambda_h(s,a)}{c_h(s,a)} \\ &= \text{A matrix game !}\end{aligned}$$

with $\mathcal{X} := \{(h,s,a) : c_h(s,a) > 0\}$ and $\Delta_{\mathcal{X}}$ = the simplex over \mathcal{X} .

$$\sum_{h,s,a} \frac{p_h^{\pi^{\text{exp}}}(s,a) \lambda_h(s,a)}{c_h(s,a)} = \mathbb{E}_{\mathcal{M}, \pi^{\text{exp}}} \left[\sum_{h,s,a} \frac{\mathbb{1}(s_h = s, a_h = a) \lambda_h(s,a)}{c_h(s,a)} \right] = \text{a value function}$$

$$\sum_{h,s,a} \frac{p_h^{\pi^{\text{exp}}}(s,a) \lambda_h(s,a)}{c_h(s,a)} = \lambda^\top (p^{\pi^{\text{exp}}} / c) = \text{linear loss of a forecaster}$$

Vanilla CovGame: pseudo-code

Algorithm 2 Vanilla CovGame

- 1: **Input:** target function c , Adversarial RL algo \mathcal{A}^Π , Online learner \mathcal{A}^λ , risk δ .
- 2: Initialize dataset $\mathcal{D}_0 \leftarrow \emptyset$
- 3: Set target set $\mathcal{X} = \{(s, a, h) \in [H] \times \mathcal{S} \times \mathcal{A} : c_h(s, a) > 0\}$
- 4: Normalize targets $\tilde{c}_h(s, a) = c_h(s, a) / c_{\min}$ ($c_{\min} := \min_{(h,s,a) \in \mathcal{X}} c_h(s, a)$)
- 5: Initialize challenger weights $\lambda_h^1(s, a) \leftarrow \mathbb{1}((h, s, a) \in \mathcal{X}) / |\mathcal{X}|$ for all h, s, a
- 6: **for** $t = 1, 2, \dots$ **do**
- 7: Define reward $R_h^t(s, a) = \mathbb{1}((h, s, a) \in \mathcal{X}) \lambda_h^t(s, a) / \tilde{c}_h(s, a)$ for all h, s, a
- 8: Feed \mathcal{A}^Π with R^t , confidence $\delta/2$ and get exploration policy π^t
- 9: Play π^t and observe trajectory $\mathcal{H}_t := \{(s_h^t, a_h^t, s_{h+1}^t)\}_{1 \leq h \leq H-1}$
- 10: Update dataset $\mathcal{D}_t \leftarrow \mathcal{D}_{t-1} \cup \mathcal{H}_t$.
- 11: Feed \mathcal{A}^λ with loss

$$\ell^t(\lambda) = \sum_{(h,s,a) \in \mathcal{X}} \lambda_h(s, a) \frac{\mathbb{1}(s_h^t = s, a_h^t = a)}{\tilde{c}_h(s, a)}$$

and get new weight vector λ^{t+1}

- 12: **If** $\forall (h, s, a), n_h(s, a; \mathcal{D}_t) \geq c_h(s, a)$: Stop and return \mathcal{D}_t

Vanilla CovGame: Sample Complexity

Theorem 2

Let $\mathcal{R}^\lambda(T)$ be an anytime bound on the regret of the online learning algorithm \mathcal{A}^λ :

$$\forall T \in \mathbb{N}^*, \sum_{t=1}^T \ell^t(\lambda^t) - \min_{\lambda \in \Delta_X} \sum_{t=1}^T \ell^t(\lambda) \leq \mathcal{R}^\lambda(T). \quad (1)$$

Let $\mathcal{R}^\Pi(T, \delta)$ be a high-probability anytime bound on the adversarial regret of \mathcal{A}^Π :

$$\mathbb{P} \left(\forall T \in \mathbb{N}^*, \sum_{t=1}^T \sup_{\pi} V_1^\pi(s_1; R^t) - \sum_{t=1}^T V_1^{\pi^t}(s_1; R^t) \leq \mathcal{R}^\Pi(T, \delta) \right) \geq 1 - \delta. \quad (2)$$

Then w.p at least $1 - \delta$, for all $T \geq 1$,

$$\min_{(h,s,a) \in \mathcal{X}} \frac{n_h^T(s,a)}{c_h(s,a)} \geq \frac{T}{\varphi^*(c)} - \frac{1}{c_{\min}} \left[\mathcal{R}^\lambda(T) + \mathcal{R}^\Pi(T, \delta/2) + \sqrt{T \log \left(\frac{4T^2}{\delta} \right)} \right]$$

Vanilla CovGame: Instanciación

Using $\mathcal{A}^\Pi = \text{UCBVI}$ (for changing rewards) and $\mathcal{A}^\lambda = \text{HEDGE}$, we have

$$\begin{aligned}\mathcal{R}^\Pi(T, \delta) &\leq 32 \log(T+1) \sqrt{SAH^2 T (\log(2SAH/\delta) + S)} \\ \mathcal{R}^\lambda(T) &\leq \sqrt{2T \log(SAH)} + \sqrt{\log(SAH)/8}\end{aligned}$$

Corollary 1

With probability at least $1 - \delta$ COVGAME instantiated with the algorithms above has sample complexity

$$\tau \leq 2\varphi^*(c) + \tilde{O}\left(\left(\frac{\varphi^*(c)}{c_{\min}}\right)^2 SH^2 A(\log(1/\delta) + S)\right)$$

Vanilla CovGame: Instanciation

Using $\mathcal{A}^\Pi = \text{UCBVI}$ (for changing rewards) and $\mathcal{A}^\lambda = \text{HEDGE}$, we have

$$\begin{aligned}\mathcal{R}^\Pi(T, \delta) &\leq 32 \log(T+1) \sqrt{SAH^2 T (\log(2SAH/\delta) + S)} \\ \mathcal{R}^\lambda(T) &\leq \sqrt{2T \log(SAH)} + \sqrt{\log(SAH)/8}\end{aligned}$$

Corollary 1

With probability at least $1 - \delta$ COVGAME instantiated with the algorithms above has sample complexity

$$\tau \leq 2\varphi^*(c) + \tilde{O}\left(\left(\frac{\varphi^*(c)}{c_{\min}}\right)^2 SH^2 A (\log(1/\delta) + S)\right)$$

When $c_h(s, a) = N \mathbb{1}((h, s, a) \in \mathcal{X})$ and $N \gg 1$, $\tau \approx 2\varphi^*(c)$

Vanilla CovGame: Proof sketch 1

Step 1: from the counts to the loss of the λ player.

$$\begin{aligned} c_{\min} \min_{(h,s,a) \in \mathcal{X}} \frac{n_h^T(s,a)}{c_h(s,a)} &= \inf_{\lambda \in \Delta_{\mathcal{X}}} \lambda \cdot (n^T / \tilde{c}) && \text{(definitions of } \tilde{c} \text{ and } \Delta_{\mathcal{X}}) \\ &= \inf_{\lambda \in \Delta_{\mathcal{X}}} \sum_{(h,s,a) \in \mathcal{X}} \lambda_h(s,a) \sum_{t=1}^T \frac{\mathbb{1}(s_h^t = s, a_h^t = a)}{\tilde{c}_h(s,a)} \\ &= \inf_{\lambda \in \Delta_{\mathcal{X}}} \sum_{t=1}^T \ell^t(\lambda) && \text{(definition of } \ell_t(\lambda)) \\ &\geq \sum_{t=1}^T \ell^t(\lambda^t) - \mathcal{R}^\lambda(T), && \text{(regret bound of } \lambda \text{ player)} \end{aligned}$$

Vanilla CovGame: Proof sketch 2

Step 2: from the loss of the λ player to the optimal value function of π player.

$$\begin{aligned} \sum_{t=1}^T \ell^t(\lambda^t) &= \sum_{t=1}^T \sum_{h,s,a} \frac{\mathbb{1}((h,s,a) \in \mathcal{X}) \lambda_h^t(s,a)}{\tilde{c}_h(s,a)} \left(\mathbb{1}(s_h^t = s, a_h^t = a) \pm p_h^{\pi^t}(s,a) \right) \\ &\hspace{20em} \text{(definition of } \ell_t(\lambda^t) \text{)} \\ &= \sum_{t=1}^T \sum_{h,s,a} p_h^{\pi^t}(s,a) R_h^t(s,a) + \sum_{t=1}^T \sum_{h,s,a} R_h^t(s,a) \left(\mathbb{1}(s_h^t = s, a_h^t = a) - p_h^{\pi^t}(s,a) \right) \\ &= \sum_{t=1}^T V_1^{\pi^t}(s_1; R^t) + M_T \hspace{10em} \text{(definition of } V_1^\pi(s_1; R) \text{ + martingale)} \\ &\geq \sup_{\pi} \sum_{t=1}^T V_1^\pi(s_1; R^t) - \mathcal{R}^\Pi(T, \delta/2) - \sqrt{T \log \left(\frac{4T^2}{\delta} \right)}, \\ &\hspace{20em} \text{(Regret of } \pi \text{ player + Azuma Hoeffding)} \end{aligned}$$

Proof sketch 3

Step 3: from the optimal value function of π player to the minimum flow.

$$\begin{aligned} \sup_{\pi} \sum_{t=1}^T V_1^{\pi}(s_1; R^t) &= \sup_{\pi} \sum_{t=1}^T \sum_{h,s,a} p_h^{\pi}(s, a) \frac{\mathbb{1}((h, s, a) \in \mathcal{X}) \lambda_h^t(s, a)}{\tilde{c}_h(s, a)} \\ &= T \sup_{\pi} \sum_{h,s,a} \left(p_h^{\pi}(s, a) \frac{\mathbb{1}((h, s, a) \in \mathcal{X})}{\tilde{c}_h(s, a)} \right) \left(\frac{\sum_{t=1}^T \lambda_h^t(s, a)}{T} \right) \\ &\geq T \sup_{\pi} \min_{(h,s,a) \in \mathcal{X}} \frac{p_h^{\pi}(s, a)}{\tilde{c}_h(s, a)} = c_{\min} \frac{T}{\varphi^*(c)}. \end{aligned}$$

Wrapping up everything, we get

$$\begin{aligned} c_{\min} \min_{(h,s,a) \in \mathcal{X}} \frac{n_h^T(s, a)}{c_h(s, a)} &\geq c_{\min} \frac{T}{\varphi^*(c)} - \mathcal{R}^{\lambda}(T) - \mathcal{R}^{\Pi}(T, \delta/2) - \sqrt{T \log \left(\frac{4T^2}{\delta} \right)} \\ \implies \min_{(h,s,a) \in \mathcal{X}} \frac{n_h^T(s, a)}{c_h(s, a)} &\geq \frac{T}{\varphi^*(c)} - \frac{1}{c_{\min}} \left[\mathcal{R}^{\lambda}(T) + \mathcal{R}^{\Pi}(T, \delta/2) + \sqrt{T \log \left(\frac{4T^2}{\delta} \right)} \right] \end{aligned}$$

Outline

- 1 Framework & definition of active coverage
- 2 Why is it worth studying ?
- 3 Sample complexity lower bound
- 4 Vanilla CovGame: a simple sometimes-optimal algorithm
- 5 Improved CovGame: a (near-)optimal algorithm**
- 6 Discussion

Improved CovGame: Idea

Recall the bound for Vanilla CovGame:

$$\tau \leq 2\varphi^*(c) + \tilde{O}\left(\left(\frac{\varphi^*(c)}{c_{\min}}\right)^2 S^2 H^2 A \log(1/\delta)\right),$$

\implies the second-order term is no longer negligible if $\frac{c_{\max}}{c_{\min}} \gg 1$!!

Improved CovGame: Idea

Recall the bound for Vanilla CovGame:

$$\tau \leq 2\varphi^*(c) + \tilde{O}\left(\left(\frac{\varphi^*(c)}{c_{\min}}\right)^2 S^2 H^2 A \log(1/\delta)\right),$$

\implies the second-order term is no longer negligible if $\frac{c_{\max}}{c_{\min}} \gg 1$!!

Idea: Cluster triplets (h, s, a) into sets

$\mathcal{Y}_j = \{(h, s, a) : c_h(s, a) \in [c_{\min}2^j, c_{\min}2^{j+1}]\}$ so that $\frac{\max_{(h,s,a) \in \mathcal{Y}_j} c_h(s,a)}{\min_{(h,s,a) \in \mathcal{Y}_j} c_h(s,a)} \leq 2$, then run Vanilla CovGame on each cluster separately.

Improved CovGame: pseudo-code

Algorithm 3 Improved CovGame

- 1: **Input:** Target function c , Adversarial RL algo \mathcal{A}^Π , Online learner \mathcal{A}^λ , risk δ .
- 2: Let $\mathcal{X}_0 := \mathcal{X}$ and $\mathcal{X}_k := \{(h, s, a) : c_h(s, a) > c_{\min} 2^k\}$ for all $k \in \mathbb{N}^*$
- 3: Initialize counts $n_h^0(s, a) = 0$ for all h, s, a
- 4: Reset \mathcal{A}^λ on $\mathcal{P}(\mathcal{X})$, set $\lambda_h^1(s, a) \leftarrow \mathbb{1}((h, s, a) \in \mathcal{X})/|\mathcal{X}|$ for all h, s, a
- 5: Initialize $k_1 \leftarrow 0$
- 6: **for** $t = 1, 2, \dots$ **do**
- 7: Get π^t from \mathcal{A}^Π given reward function λ^t and confidence $1 - \delta/2$
- 8: Generate a trajectory $\{(s_h^t, a_h^t)\}_{h \in [H]}$ using policy π^t and update counts n^t
- 9: **if** $n_h^t(s, a) \geq c_h(s, a)$ for all h, s, a **then** stop and return all sampled trajectories
- 10: Update $k_{t+1} \leftarrow \max\{j \in \mathbb{N} : n_h^t(s, a) \geq c_h(s, a) \forall (h, s, a) \in \mathcal{X} \setminus \mathcal{X}_j\}$
- 11: **if** $k_{t+1} \neq k_t$ **then**
- 12: Reset \mathcal{A}^λ on $\mathcal{P}(\mathcal{X}_{k_{t+1}})$, set $\lambda_h^{t+1}(s, a) \leftarrow \mathbb{1}((h, s, a) \in \mathcal{X}_{k_{t+1}})/|\mathcal{X}_{k_{t+1}}|$
- 13: **else**
- 14: Feed \mathcal{A}^λ with loss $\ell^t(\lambda) = \sum_{(h,s,a) \in \mathcal{X}_{k_t}} \lambda_h(s, a) \mathbb{1}(s_h^t = s, a_h^t = a)$, get weight λ^{t+1}

Corollary 2 (Under the same conditions of Corollary 1)

With probability at least $1 - \delta$, the stopping time of Improved CovGame with HEDGE and UCBVI can be bounded by

$$\tau \leq 8m\varphi^*(c) + \tilde{O}(m^2\varphi^*(\mathbb{1}_X)^2SAH^2(\log(1/\delta) + S)),$$

where $m = \lceil \log_2(c_{\max}/c_{\min}) \rceil \vee 1$.

Corollary 2 (Under the same conditions of Corollary 1)

With probability at least $1 - \delta$, the stopping time of Improved CovGame with HEDGE and UCBVI can be bounded by

$$\tau \leq 8m\varphi^*(c) + \tilde{O}(m^2\varphi^*(1_{\mathcal{X}})^2SAH^2(\log(1/\delta) + S)),$$

where $m = \lceil \log_2(c_{\max}/c_{\min}) \rceil \vee 1$.

\implies The scaling in c_{\max}/c_{\min} is only logarithmic.

Outline

- 1 Framework & definition of active coverage
- 2 Why is it worth studying ?
- 3 Sample complexity lower bound
- 4 Vanilla CovGame: a simple sometimes-optimal algorithm
- 5 Improved CovGame: a (near-)optimal algorithm
- 6 Discussion

Previous approaches

- (Tarbouriech et al. 2021) GOSPRL
- (Wagenmaker et al. 2021) Learn2Explore

Previous approaches

- (Tarbouriech et al. 2021) GOSPRL
- (Wagenmaker et al. 2021) Learn2Explore

Complexity:

$$\sum_{h,s,a} \frac{c_h(s,a)}{\sup_{\pi^{exp}} p_h^{\pi^{exp}}(s,a)} \geq \varphi^*(c)$$

Conclusion and perspectives

- 1 Active coverage = how collect $c_h(s, a)$ observations from each (h, s, a) of an MDP.
- 2 The minimum number of episode one needs is
$$\varphi^*(c) = \inf_{\pi^{exp} \in \Pi^S} \max_{(s,a,h) \in \mathcal{X}} \frac{c_h(s,a)}{p_h^{\pi^{exp}}(s,a)}$$
, a quantity that depends on the MDP.
- 3 A simple game-based strategy can (nearly) match this lower bound.

Thanks !