

Adaptive Sampling for Best Policy Identification in MDPs

Aymen Al Marjani¹
joint work with Alexandre Proutiere²

¹ENS de Lyon

²KTH Royal Institute of Technology

38th International Conference on Machine Learning

- 1 Introduction
- 2 Information-Theoretical Lower Bound
- 3 Upper bound of the characteristic time
- 4 Algorithm Design
- 5 Experiments
- 6 Conclusion

How many samples does it take to learn an optimal policy in RL ?

Infinite horizon discounted MDPs

$$\phi = \langle \mathcal{S}, \mathcal{A}, p_\phi, q_\phi, \gamma \rangle$$

- ① \mathcal{S}, \mathcal{A} : **Finite** state and action spaces.

Infinite horizon discounted MDPs

$$\phi = \langle \mathcal{S}, \mathcal{A}, p_\phi, q_\phi, \gamma \rangle$$

- 1 \mathcal{S}, \mathcal{A} : **Finite** state and action spaces.
- 2 After choosing action a at state s the agent:
 - receives reward $R(s, a) \sim q_\phi(\cdot | s, a)$ and mean $r(s, a) \triangleq \mathbb{E}_{q(\cdot | s, a)}[R(s, a)]$.
 - makes transition to $s' \sim p_\phi(\cdot | s, a)$.

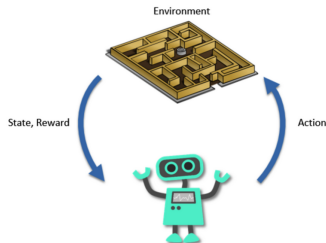


Figure: src:packtpub

Infinite horizon discounted MDPs

$$\phi = \langle \mathcal{S}, \mathcal{A}, p_\phi, q_\phi, \gamma \rangle$$

- 1 \mathcal{S}, \mathcal{A} : **Finite** state and action spaces.
- 2 After choosing action a at state s the agent:
 - receives reward $R(s, a) \sim q_\phi(\cdot | s, a)$ and mean $r(s, a) \triangleq \mathbb{E}_{q_\phi(\cdot | s, a)}[R(s, a)]$.
 - makes transition to $s' \sim p_\phi(\cdot | s, a)$.
 - For simplicity, we assume q with support in $[0, 1]$.

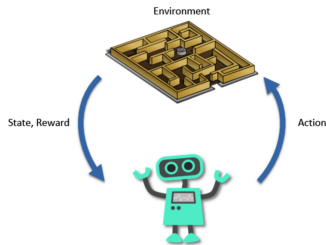


Figure: src:packtpub

Best Policy Identification

$$\phi = \langle \mathcal{S}, \mathcal{A}, p_\phi, q_\phi, \gamma \rangle$$

- $\gamma \in [0, 1)$ is the discount factor.

Best Policy Identification

$\phi = \langle \mathcal{S}, \mathcal{A}, p_\phi, q_\phi, \gamma \rangle$

- $\gamma \in [0, 1)$ is the discount factor.
- Identify a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ maximizing the total discounted reward:

$$\pi_\phi^* \in \arg \max_{\pi} V_\phi^\pi(s) = \mathbb{E}_\phi \left[\sum_{t=0}^{\infty} \gamma^t R(s_t^\pi, \pi(s_t^\pi)) \mid s_0 = s \right]$$

- **δ -PC algorithm:** $\mathbb{P}_\phi(\widehat{\pi}_\tau^* \neq \pi^*) \leq \delta$.

Best Policy Identification

$\phi = \langle \mathcal{S}, \mathcal{A}, p_\phi, q_\phi, \gamma \rangle$

- $\gamma \in [0, 1)$ is the discount factor.
- Identify a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ maximizing the total discounted reward:

$$\pi_\phi^* \in \arg \max_{\pi} V_\phi^\pi(s) = \mathbb{E}_\phi \left[\sum_{t=0}^{\infty} \gamma^t R(s_t^\pi, \pi(s_t^\pi)) \mid s_0 = s \right]$$

- **δ -PC algorithm:** $\mathbb{P}_\phi(\hat{\pi}_\tau^* \neq \pi^*) \leq \delta$.
- Identify π^* **using minimum number of samples!**

- **Assumption 1:** $\pi^* \triangleq \pi_{\phi}^*$ is unique.

- **Assumption 1:** $\pi^* \triangleq \pi_\phi^*$ is unique.
- **Generative Model:** The agent has access to a simulator. At round t , she agent can query a sample *any* pair (s_t, a_t) . She then observes $(R_t, s'_t) \sim q_\phi(\cdot | s_t, a_t) \otimes p_\phi(\cdot | s_t, a_t)$. Next, she can choose *any* other pair (s_{t+1}, a_{t+1}) *independently of her previous state*.

Learning: be specific!

Two kinds of guarantees:

- **Minimax** over a set of MDPs Φ :

$$\inf_{\mathbb{A}: \delta\text{-PC}} \sup_{\phi \in \Phi} \mathbb{E}_{\phi, \mathbb{A}}[\tau_\delta]$$

- **Instance-specific:** For a given ϕ :

$$\inf_{\mathbb{A}: \delta\text{-PC}} \mathbb{E}_{\phi, \mathbb{A}}[\tau_\delta]$$

Learning: be specific!

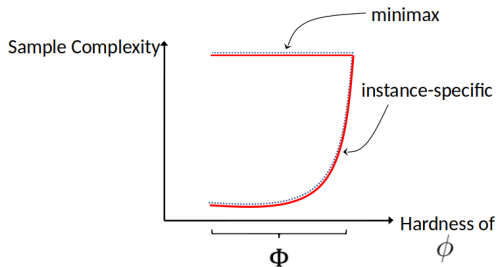
Two kinds of guarantees:

- **Minimax** over a set of MDPs Φ :

$$\inf_{\mathbb{A}: \delta\text{-PC}} \sup_{\phi \in \Phi} \mathbb{E}_{\phi, \mathbb{A}}[\tau_{\delta}]$$

- **Instance-specific:** For a given ϕ :

$$\inf_{\mathbb{A}: \delta\text{-PC}} \mathbb{E}_{\phi, \mathbb{A}}[\tau_{\delta}]$$



Learning: be specific!

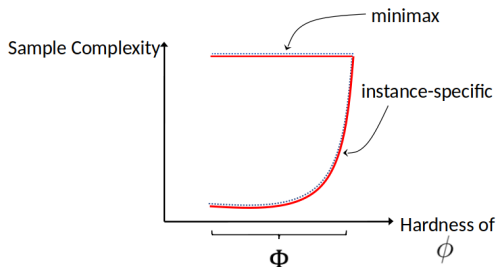
Two kinds of guarantees:

- **Minimax** over a set of MDPs Φ :

$$\inf_{\mathbb{A}: \delta\text{-PC}} \sup_{\phi \in \Phi} \mathbb{E}_{\phi, \mathbb{A}}[\tau_{\delta}]$$

- **Instance-specific:** For a given ϕ :

$$\inf_{\mathbb{A}: \delta\text{-PC}} \mathbb{E}_{\phi, \mathbb{A}}[\tau_{\delta}]$$



- We seek algorithms that can adapt to the hardness of the instance.

Define:

- The set of alternative MDPs $\text{Alt}(\phi) = \{\psi : \pi^* \text{ is not optimal in } \psi\}$.
- Σ the simplex of \mathbb{R}^{SA} .
- $\text{KL}_{\phi|\psi}(s, a) = \text{KL}(q_{\phi}(s, a), q_{\psi}(s, a)) + \text{KL}(p_{\phi}(s, a), p_{\psi}(s, a))$

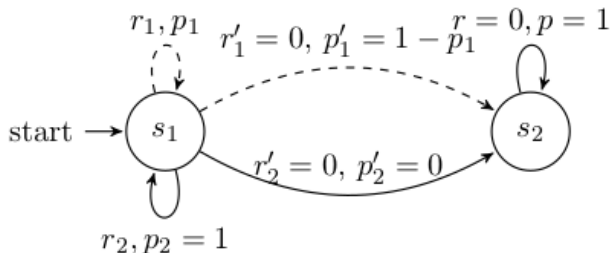
Proposition 1

The sample complexity of any δ -PC algorithm satisfies: for any ϕ with a unique optimal policy,

$$\mathbb{E}_{\phi}[\tau] \geq T^*(\phi) \log(1/2.4\delta),$$

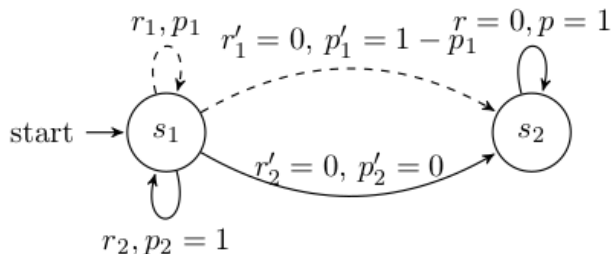
$$\text{where } T^*(\phi)^{-1} = \sup_{\omega \in \Sigma} \inf_{\psi \in \text{Alt}(\phi)} \sum_{s,a} \omega_{sa} \text{KL}_{\phi|\psi}(s, a). \quad (1)$$

IT Lower bound: Hard to solve !



- $\text{Alt}(\phi)$ and $\text{Alt}_{s_1 a_1}(\phi)$ are not convex.

IT Lower bound: Hard to solve !



- $\text{Alt}(\phi)$ and $\text{Alt}_{s_1 a_1}(\phi)$ are not convex.
- \implies The sub-problem $\inf_{\psi \in \text{Alt}(\phi)} \sum_{s,a} \omega_{sa} \text{KL}_{\phi|\psi}(s, a)$ is non-convex.

IT Lower bound: MDP vs MAB

	MAB	MDP
Parameters	$\mu_1 > \dots \geq \mu_K$	$(r(s, a), p(s, a))_{s,a}$
Objective	Identify $a^* = \arg \max_{a \in [K]} \mu_a$	Identify $\pi^* = \arg \max_{\pi} (I - \gamma P_{\pi})^{-1} r_{\pi}$
Alternative instances	$\bigcup_{a \neq 1} \{\lambda : \lambda_a > \lambda_1\}$ union of convex sets	$\bigcup_{s, a \neq \pi^*(s)} \{\psi : Q_{\psi}^{\pi^*}(s, a) > V_{\psi}^{\pi^*}(s)\}$ Not union of convex sets
IT lower bound	Easy to solve	Hard to solve

Upper bound: Idea

Define the characteristic time: $T(\phi, \omega)^{-1} \triangleq \inf_{\psi \in \text{Alt}(\phi)} \sum_{\mathbf{s}, \mathbf{a}} \omega_{\mathbf{s}\mathbf{a}} \text{KL}_{\phi|\psi}(\mathbf{s}, \mathbf{a})$.

Upper bound: Idea

Define the characteristic time: $T(\phi, \omega)^{-1} \triangleq \inf_{\psi \in \text{Alt}(\phi)} \sum_{s,a} \omega_{sa} \text{KL}_{\phi|\psi}(s, a)$.

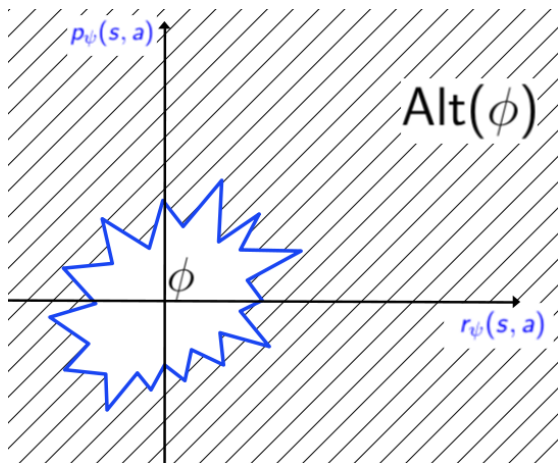


Figure: $\text{Alt}(\phi)$: Non-convex boundary

Upper bound: Idea

Define the characteristic time: $T(\phi, \omega)^{-1} \triangleq \inf_{\psi \in \text{Alt}(\phi)} \sum_{s,a} \omega_{sa} \text{KL}_{\phi|\psi}(s, a)$.

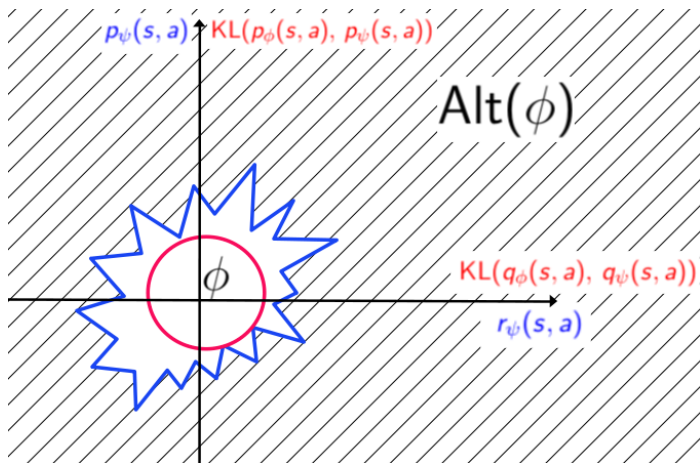


Figure: KL Ball

Upper bound: Instance-specific quantities

Define:

- The sub-optimality gap: $\Delta_{sa} = V_{\phi}^*(s) - Q_{\phi}^*(s, a)$.
- The minimum gap $\Delta_{\min} = \min_{s, a \neq \pi^*(s)} \Delta_{sa}$.
- The variance of the value function $\text{Var}_{(s,a)}[V_{\phi}^*] = \mathbb{V}_{s' \sim p_{\phi}(\cdot|s,a)}[V_{\phi}^*(s)]$.
- The span of the value function $\text{sp}[V_{\phi}^*] = \max_s V_{\phi}^*(s) - \min_s V_{\phi}^*(s)$.

Upper bound of the characteristic time

Theorem 1 (Upper bound of minimal sample complexity)

For all vectors ω in the simplex:

$$T(\phi, \omega) \leq U(\phi, \omega) \triangleq \max_{s, a \neq \pi^*(s)} \frac{T_1(s, a; \phi) + T_2(s, a; \phi)}{\omega_{sa}} + \frac{T_3(\phi) + T_4(\phi)}{\min_s \omega_{s, \pi^*(s)}},$$

where

$$\begin{cases} T_1(s, a; \phi) = \frac{2}{\Delta_{sa}^2}, \\ T_2(s, a; \phi) = \max \left(\frac{16 \text{Var}_{(s,a)}[V_\phi^*]}{\Delta_{sa}^2}, \frac{6 \text{sp}[V_\phi^*]^{4/3}}{\Delta_{sa}^{4/3}} \right), \\ T_3(\phi) = \frac{2}{[\Delta_{\min}(\phi)(1-\gamma)]^2}, \\ T_4(\phi) \leq \frac{27}{\Delta_{\min}(\phi)^2(1-\gamma)^3} = \mathcal{O} \left(\frac{\text{Minimax lower bound}}{SA} \right) \end{cases}$$

- The optimal weights minimizing the upper-bound program:

$$\bar{\omega}(\phi) = \arg \inf_{\omega \in \Sigma} \max_{(s,a): a \neq \pi^*(s)} \frac{T_1(s, a; \phi) + T_2(s, a; \phi)}{\omega_{sa}} + \frac{T_3(\phi) + T_4(\phi)}{\min_s \omega_{s, \pi^*(s)}}$$

are easy to compute !

- The optimal weights minimizing the upper-bound program:

$$\bar{\omega}(\phi) = \arg \inf_{\omega \in \Sigma} \max_{(s,a): a \neq \pi^*(s)} \frac{T_1(s, a; \phi) + T_2(s, a; \phi)}{\omega_{sa}} + \frac{T_3(\phi) + T_4(\phi)}{\min_s \omega_{s, \pi^*(s)}}$$

are **easy to compute** !

- Ensures that $\mathbb{P}_\phi \left(\forall (s, a) \in \mathcal{S} \times \mathcal{A}, \lim_{t \rightarrow \infty} \frac{N_{sa}(t)}{t} = \bar{\omega}_{s,a}(\phi) \right) = 1$.

KLB-TS: stopping rule

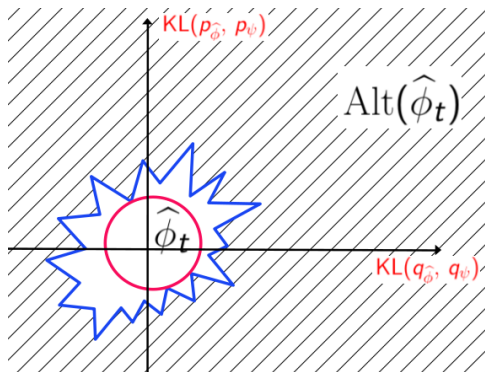


Figure: KL-Ball Stopping rule

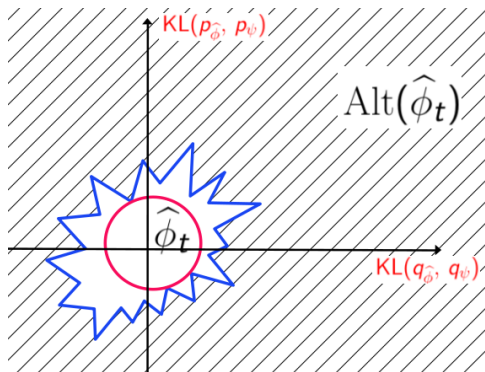


Figure: KL-Ball Stopping rule

- We ensure that ϕ falls within the KL-ball with probability $1 - \delta$, using PAC bounds on the KL divergence..

Theorem 3

KLB-TS has a sample complexity τ_δ satisfying:

for all $\delta \in (0, 1)$, $\mathbb{E}_\phi[\tau_\delta]$ is finite and $\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_\phi[\tau_\delta]}{\log(1/\delta)} \leq 4U(\phi)$, where:

$$\begin{aligned} U(\phi) &\triangleq \inf_{\omega} U(\phi, \omega) \\ &= \mathcal{O}\left(S \min\left(\frac{\text{Var}_{\max}^*[V_\phi^*]}{\Delta_{\min}^2 (1-\gamma)^2}, \frac{1}{\Delta_{\min}^2 (1-\gamma)^3}\right)\right) \\ &\quad + \sum_{s, a \neq \pi^*(s)} \frac{1 + \text{Var}_{(s,a)}[V_\phi^*]}{\Delta_{sa}^2} \end{aligned}$$

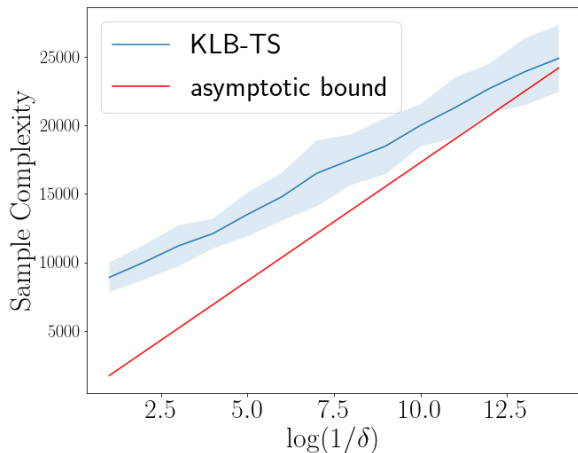


Figure: Asymptotic bound: $S=A=2$, $\gamma = 0.5$.

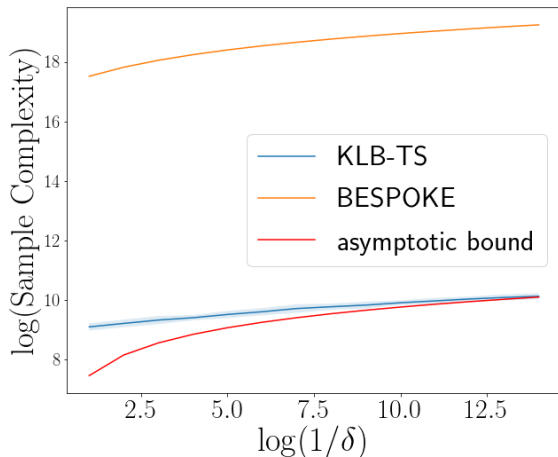


Figure: KLB-TS vs. BESPOKE. $S=A=2$, $\gamma = 0.5$.

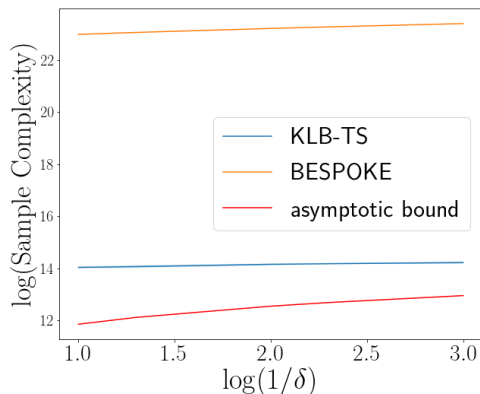


Figure: KLB-TS vs. BESPOKE. $S = 5, A = 10, \gamma = 0.7$.

- 1 Contrary to MAB, IT lower bound is hard to solve for MDPs.
- 2 We can derive problem-specific surrogates which :
 - Are *explicit*, depending on functionals of the MDP.
 - Have a corresponding allocation that is easy to compute.
- 3 Can be used to devise (Asymptotically) Matching algorithm.
- 4 First step towards understanding problem-specific ε -optimal policy identification.

Thanks !