

# Adaptive Sampling for Best Policy Identification in MDPs

Aymen Al Marjani<sup>1</sup>  
joint work with Alexandre Proutiere<sup>2</sup>

<sup>1</sup>ENS de Lyon

<sup>2</sup>KTH Royal Institute of Technology

RL Theory seminars

May 11th, 2021

# Outline

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

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

# Infinite horizon MDPs

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

# Infinite horizon 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 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  $s' \sim p_\phi(\cdot | s, a)$ .

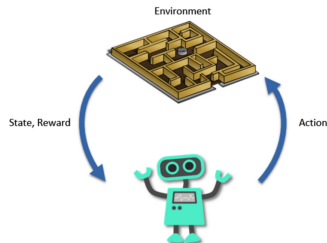


Figure: src:packtpub

# Infinite horizon 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  $[0, 1]$  .

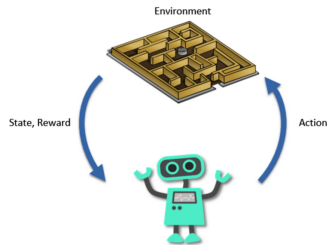


Figure: src:packtpub

# 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:

# 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:

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

# 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:

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

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

- **Forward model:** The agent can only follow trajectories:  $(s_0, a_0, R_0, s_1, a_1 \dots, )$  where  $s_{t+1} \sim p_\phi(\cdot | s_t, a_t)$ .

- **Forward model:** The agent can only follow trajectories:  $(s_0, a_0, R_0, s_1, a_1 \dots, )$  where  $s_{t+1} \sim p_\phi(\cdot | s_t, a_t)$ .
- **Generative model:** At round  $t$ , the agent can 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*.

- **Forward model:** The agent can only follow trajectories:  $(s_0, a_0, R_0, s_1, a_1 \dots, )$  where  $s_{t+1} \sim p_\phi(\cdot | s_t, a_t)$ .
- **Generative model:** At round  $t$ , the agent can 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*.

In this talk, we focus on the Generative model.

- **Sampling rule:** How to select next pair to sample depending on past observations:  $(s_{t+1}, a_{t+1})$  is  $\mathcal{F}_t \triangleq \sigma((s_j, a_j, R_j, s'_j)_{1 \leq j \leq t})$  measurable.

- **Sampling rule:** How to select next pair to sample depending on past observations:  $(s_{t+1}, a_{t+1})$  is  $\mathcal{F}_t \triangleq \sigma((s_j, a_j, R_j, s'_j)_{1 \leq j \leq t})$  measurable.
- **Stopping rule:** The algorithm stops sampling after collecting  $\tau$  samples and returns  $\hat{\pi}_\tau^*$ .  $\tau$  is a stopping time w.r.t. the filtration  $(\mathcal{F}_t)_{t \geq 1}$ .



- **Sampling rule:** How to select next pair to sample depending on past observations:  $(s_{t+1}, a_{t+1})$  is  $\mathcal{F}_t \triangleq \sigma((s_j, a_j, R_j, s'_j)_{1 \leq j \leq t})$  measurable.
- **Stopping rule:** The algorithm stops sampling after collecting  $\tau$  samples and returns  $\hat{\pi}_\tau^*$ .  $\tau$  is a stopping time w.r.t. the filtration  $(\mathcal{F}_t)_{t \geq 1}$ .
- **$\delta$ -PC algorithm:**  $\mathbb{P}_\phi(\hat{\pi}_\tau^* \neq \pi^*) \leq \delta$ .

- **Sampling rule:** How to select next pair to sample depending on past observations:  $(s_{t+1}, a_{t+1})$  is  $\mathcal{F}_t \triangleq \sigma((s_j, a_j, R_j, s'_j)_{1 \leq j \leq t})$  measurable.
- **Stopping rule:** The algorithm stops sampling after collecting  $\tau$  samples and returns  $\hat{\pi}_\tau^*$ .  $\tau$  is a stopping time w.r.t. the filtration  $(\mathcal{F}_t)_{t \geq 1}$ .
- **$\delta$ -PC algorithm:**  $\mathbb{P}_\phi(\hat{\pi}_\tau^* \neq \pi^*) \leq \delta$ .
- Identify  $\pi^*$  as fast as possible!

- **Sampling rule:** How to select next pair to sample depending on past observations:  $(s_{t+1}, a_{t+1})$  is  $\mathcal{F}_t \triangleq \sigma((s_j, a_j, R_j, s'_j)_{1 \leq j \leq t})$  measurable.
- **Stopping rule:** The algorithm stops sampling after collecting  $\tau$  samples and returns  $\hat{\pi}_\tau^*$ .  $\tau$  is a stopping time w.r.t. the filtration  $(\mathcal{F}_t)_{t \geq 1}$ .
- **$\delta$ -PC algorithm:**  $\mathbb{P}_\phi(\hat{\pi}_\tau^* \neq \pi^*) \leq \delta$ .
- Identify  $\pi^*$  **as fast as possible!**

$\implies$  Algorithm with minimal *sample complexity*  $\tau$

# Learning: be specific!

Two kinds of guarantees:

# Learning: be specific!

Two kinds of guarantees:

- **Minimax** over a set of MDPs  $\Phi$ :

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

# Learning: be specific!

Two kinds of guarantees:

- **Minimax** over a set of MDPs  $\Phi$ :

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

- Minimax lower bounds often come from pathological examples. Real world scenarios are not that hard (unless in adversarial settings).

# Learning: be specific!

Two kinds of guarantees:

- **Minimax** over a set of MDPs  $\Phi$ :

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

- Minimax lower bounds often come from pathological examples. Real world scenarios are not that hard (unless in adversarial settings).
- Algorithms that sample state-actions uniformly at random are sufficient to be minimax optimal !

# Learning: be specific!

Two kinds of guarantees:

- **Minimax** over a set of MDPs  $\Phi$ :

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

- **Instance-specific:** For a given  $\phi$ :

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



# Learning: be specific!

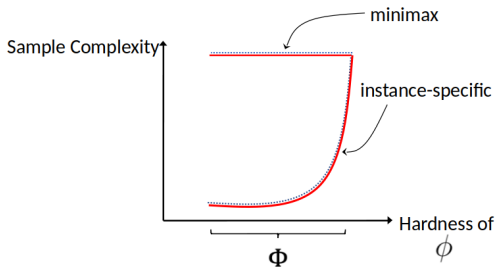
Two kinds of guarantees:

- **Minimax** over a set of MDPs  $\Phi$ :

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

- **Instance-specific:** For a given  $\phi$ :

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



# Learning: be specific!

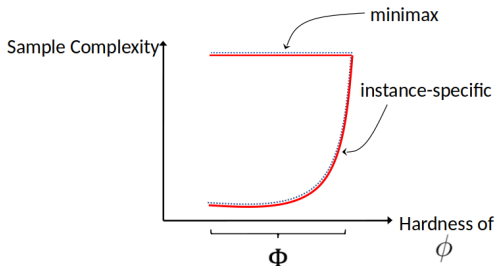
Two kinds of guarantees:

- **Minimax** over a set of MDPs  $\Phi$ :

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

- **Instance-specific:** For a given  $\phi$ :

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



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

### ① **Minimax Approach:**

- Introduced by [Kearns and Singh, 1999].

## ① **Minimax Approach:**

- Introduced by [Kearns and Singh, 1999].
- Lower bound by [Azar et al., 2013]:  $\mathcal{O}\left(\frac{SA \log(SA/\delta)}{\varepsilon^2(1-\gamma)^3}\right)$  samples to get  $\varepsilon$ -optimal policy.

## ① Minimax Approach:

- Introduced by [Kearns and Singh, 1999].
- Lower bound by [Azar et al., 2013]:  $\mathcal{O}\left(\frac{SA \log(SA/\delta)}{\varepsilon^2(1-\gamma)^3}\right)$  samples to get  $\varepsilon$ -optimal policy.
- variety of minimax-optimal algorithms (model-free and model based): [Azar et al., 2013, Sidford et al., 2018, Agarwal et al., 2020, Li et al., 2020]

## ① Minimax Approach:

- Introduced by [Kearns and Singh, 1999].
- Lower bound by [Azar et al., 2013]:  $\mathcal{O}\left(\frac{SA \log(SA/\delta)}{\varepsilon^2(1-\gamma)^3}\right)$  samples to get  $\varepsilon$ -optimal policy.
- variety of minimax-optimal algorithms (model-free and model based): [Azar et al., 2013, Sidford et al., 2018, Agarwal et al., 2020, Li et al., 2020]

## ② Problem-specific approach:

- Multi-armed Bandit:
  - [Even-Dar and Mansour, 2003, Gabillon et al., 2012, Kalyanakrishnan et al., 2012] Bounds depending on the gaps.

## ① **Minimax Approach:**

- Introduced by [Kearns and Singh, 1999].
- Lower bound by [Azar et al., 2013]:  $\mathcal{O}\left(\frac{SA \log(SA/\delta)}{\varepsilon^2(1-\gamma)^3}\right)$  samples to get  $\varepsilon$ -optimal policy.
- variety of minimax-optimal algorithms (model-free and model based): [Azar et al., 2013, Sidford et al., 2018, Agarwal et al., 2020, Li et al., 2020]

## ② **Problem-specific approach:**

- Multi-armed Bandit:
  - [Even-Dar and Mansour, 2003, Gabillon et al., 2012, Kalyanakrishnan et al., 2012] Bounds depending on the gaps.
  - [Garivier and Kaufmann, 2016] complete characterization for exponential family.

## 1 Minimax Approach:

- Introduced by [Kearns and Singh, 1999].
- Lower bound by [Azar et al., 2013]:  $\mathcal{O}\left(\frac{SA \log(SA/\delta)}{\varepsilon^2(1-\gamma)^3}\right)$  samples to get  $\varepsilon$ -optimal policy.
- variety of minimax-optimal algorithms (model-free and model based): [Azar et al., 2013, Sidford et al., 2018, Agarwal et al., 2020, Li et al., 2020]

## 2 Problem-specific approach:

- Multi-armed Bandit:
  - [Even-Dar and Mansour, 2003, Gabillon et al., 2012, Kalyanakrishnan et al., 2012] Bounds depending on the gaps.
  - [Garivier and Kaufmann, 2016] complete characterization for exponential family.
- MDPs: [Zanette et al., 2019] proposed BESPOKE, first algorithm with problem-specific guarantees.



## Related work: BESPOKE

- ① Principle: Minimize a weighted sum of confidence intervals over state-action pairs.
- ② Advantages:

## Related work: BESPOKE

- ① Principle: Minimize a weighted sum of confidence intervals over state-action pairs.
- ② Advantages:
  - Provides a clear stopping rule to find an  $\varepsilon$ -optimal policy.

# Related work: BESPOKE

- 1 Principle: Minimize a weighted sum of confidence intervals over state-action pairs.
- 2 Advantages:
  - Provides a clear stopping rule to find an  $\varepsilon$ -optimal policy.
  - First problem-specific bound, w.h.p:

$$\tau_\delta \leq \tilde{\mathcal{O}} \left( \sum_{s \in \mathcal{S}} \min \left( \frac{1}{(1-\gamma)^3 \Delta_{\min}^2}, \frac{\text{Var}_{(s, \pi^*(s))}[R] + \gamma^2 \text{Var}_{p(s, \pi^*(s))}[V_\phi^*]}{\Delta_{\min}^2} \right) + \sum_{s, a \neq \pi^*(s)} \frac{\text{Var}[R(s, a)] + \gamma^2 \text{Var}_{p(s, a)}[V_\phi^*]}{\Delta_{sa}^2} + \frac{S^2 A}{(1-\gamma)^2} \right).$$

# Related work: BESPOKE

- 1 Principle: Minimize a weighted sum of confidence intervals over state-action pairs.
- 2 Advantages:
  - Provides a clear stopping rule to find an  $\varepsilon$ -optimal policy.
  - First problem-specific bound, w.h.p:

$$\tau_\delta \leq \tilde{O} \left( \sum_{s \in \mathcal{S}} \min \left( \frac{1}{(1-\gamma)^3 \Delta_{\min}^2}, \frac{\text{Var}_{(s, \pi^*(s))}[R] + \gamma^2 \text{Var}_{p(s, \pi^*(s))}[V_\phi^*]}{\Delta_{\min}^2} \right) + \sum_{s, a \neq \pi^*(s)} \frac{\text{Var}[R(s, a)] + \gamma^2 \text{Var}_{p(s, a)}[V_\phi^*]}{\Delta_{sa}^2} + \frac{S^2 A}{(1-\gamma)^2} \right).$$

- 3 Drawbacks:
  - Solves a convex problem at every step.

- 1 Principle: Minimize a weighted sum of confidence intervals over state-action pairs.
- 2 Advantages:
  - Provides a clear stopping rule to find an  $\varepsilon$ -optimal policy.
  - First problem-specific bound, w.h.p:

$$\tau_\delta \leq \tilde{O} \left( \sum_{s \in \mathcal{S}} \min \left( \frac{1}{(1-\gamma)^3 \Delta_{\min}^2}, \frac{\text{Var}_{(s, \pi^*(s))}[R] + \gamma^2 \text{Var}_{p(s, \pi^*(s))}[V_\phi^*]}{\Delta_{\min}^2} \right) + \sum_{s, a \neq \pi^*(s)} \frac{\text{Var}[R(s, a)] + \gamma^2 \text{Var}_{p(s, a)}[V_\phi^*]}{\Delta_{sa}^2} + \frac{S^2 A}{(1-\gamma)^2} \right).$$

- 3 Drawbacks:
  - Solves a convex problem at every step.
  - Large burn-in phase:  $\Omega\left(\frac{S^2 A \log(1/\delta)}{(1-\gamma)^2}\right)$ .

# Information-Theoretical lower bound

Define:

- The set of alternative MDPs  $\text{Alt}(\phi) = \{\psi : \pi^* \text{ is not optimal in } \psi\}$ .

# Information-Theoretical lower bound

Define:

- The set of alternative MDPs  $\text{Alt}(\phi) = \{\psi : \pi^* \text{ is not optimal in } \psi\}$ .
- $\Sigma$  the simplex of  $\mathbb{R}^{SA}$ .



Define:

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

Define:

- The set of alternative MDPs  $\text{Alt}(\phi) = \{\psi : \pi^* \text{ is not optimal in } \psi\}$ .
- $\Sigma$  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  $\delta$ -PC algorithm satisfies: for any  $\phi$  with a unique optimal policy,

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

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

# Solving the lower bound program?

## Solving the lower bound program?

- By definition:  $\text{Alt}(\phi) = \{\psi : \exists(\mathbf{s}, \pi) \in \mathcal{S} \times \Pi, V_{\psi}^{\pi}(\mathbf{s}) > V_{\psi}^{\pi^*}(\mathbf{s})\}$ .

# Solving the lower bound program?

- By definition:  $\text{Alt}(\phi) = \{\psi : \exists(s, \pi) \in \mathcal{S} \times \Pi, V_{\psi}^{\pi}(s) > V_{\psi}^{\pi^*}(s)\}$ .
- Involves many parameters of  $\psi$ :

$$\left( r(x, \pi(x)), p(x, \pi(x)), r(x, \pi^*(x)), p(x, \pi^*(x)) \right)_{x \in \mathcal{S}}.$$

# Solving the lower bound program?

- By definition:  $\text{Alt}(\phi) = \{\psi : \exists (s, \pi) \in \mathcal{S} \times \Pi, V_{\psi}^{\pi}(s) > V_{\psi}^{\pi^*}(s)\}$ .
- Involves many parameters of  $\psi$ :

$$\left( r(x, \pi(x)), p(x, \pi(x)), r(x, \pi^*(x)), p(x, \pi^*(x)) \right)_{x \in \mathcal{S}}.$$

$\Rightarrow$  We need further simplification.

# Solving the lower bound program?

- By definition:  $\text{Alt}(\phi) = \{\psi : \exists (s, \pi) \in \mathcal{S} \times \Pi, V_{\psi}^{\pi}(s) > V_{\psi}^{\pi^*}(s)\}$ .
- Involves many parameters of  $\psi$ :

$$\left( r(x, \pi(x)), p(x, \pi(x)), r(x, \pi^*(x)), p(x, \pi^*(x)) \right)_{x \in \mathcal{S}}.$$

$\Rightarrow$  We need further simplification.

## Lemma 2

The set of alternative MDPs can be decomposed as follows:

$$\text{Alt}(\phi) = \bigcup_{(s,a): a \neq \pi^*(s)} \{\psi : Q_{\psi}^{\pi^*}(s, a) > V_{\psi}^{\pi^*}(s)\}. \quad (2)$$

# Solving the lower bound program?

- By definition:  $\text{Alt}(\phi) = \{\psi : \exists (s, \pi) \in \mathcal{S} \times \Pi, V_{\psi}^{\pi}(s) > V_{\psi}^{\pi^*}(s)\}$ .
- Involves many parameters of  $\psi$ :

$$\left( r(x, \pi(x)), p(x, \pi(x)), r(x, \pi^*(x)), p(x, \pi^*(x)) \right)_{x \in \mathcal{S}}.$$

⇒ We need further simplification.

## Lemma 2

The set of alternative MDPs can be decomposed as follows:

$$\text{Alt}(\phi) = \bigcup_{(s,a): a \neq \pi^*(s)} \{\psi : Q_{\psi}^{\pi^*}(s, a) > V_{\psi}^{\pi^*}(s)\}. \quad (2)$$

- In contrast with  $Q_{\phi}^{\pi^*}(s, a) < V_{\phi}^{\pi^*}(s)$ , for  $a \neq \pi^*(s)$ .



# Solving the lower bound program?

- By definition:  $\text{Alt}(\phi) = \{\psi : \exists (s, \pi) \in \mathcal{S} \times \Pi, V_{\psi}^{\pi}(s) > V_{\psi}^{\pi^*}(s)\}$ .
- Involves many parameters of  $\psi$ :

$$\left( r(x, \pi(x)), p(x, \pi(x)), r(x, \pi^*(x)), p(x, \pi^*(x)) \right)_{x \in \mathcal{S}}.$$

$\Rightarrow$  We need further simplification.

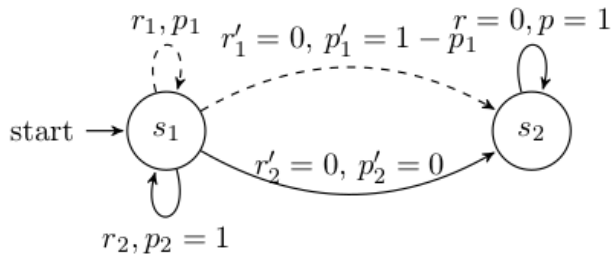
## Lemma 2

The set of alternative MDPs can be decomposed as follows:

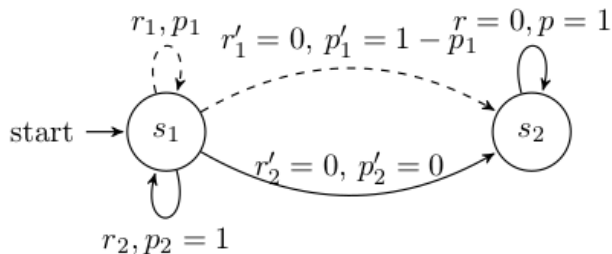
$$\text{Alt}(\phi) = \bigcup_{(s,a): a \neq \pi^*(s)} \{\psi : Q_{\psi}^{\pi^*}(s, a) > V_{\psi}^{\pi^*}(s)\}. \quad (2)$$

- In contrast with  $Q_{\phi}^{\pi^*}(s, a) < V_{\phi}^{\pi^*}(s)$ , for  $a \neq \pi^*(s)$ .
- Only involves  $(r(s, a), p(s, a))$  and  $(r(x, \pi^*(x)), p(x, \pi^*(x)))_{x \in \mathcal{S}}$  in  $\psi$ .

# IT Lower bound: Intractable!

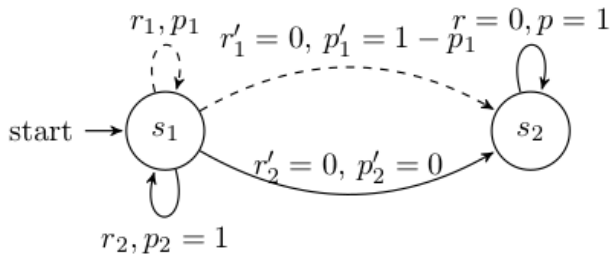


# IT Lower bound: Intractable!



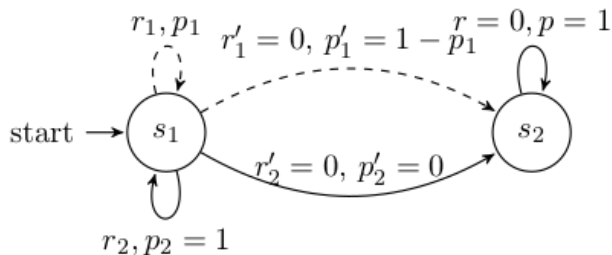
- $Q(s_1, a_i) = \frac{r_i}{1 - \gamma p_i}$ ,  $i = 1, 2$ .

# IT Lower bound: Intractable!



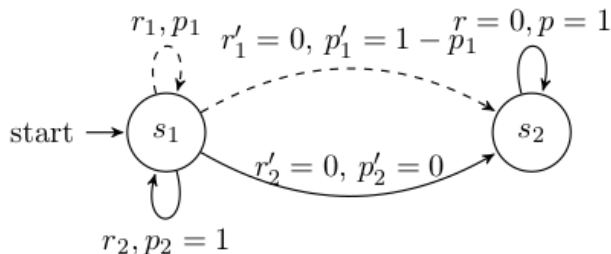
- $Q(s_1, a_i) = \frac{r_i}{1 - \gamma p_i}$ ,  $i = 1, 2$ .
- Can easily construct  $\psi$  and  $\bar{\psi}$  such that:

# IT Lower bound: Intractable!



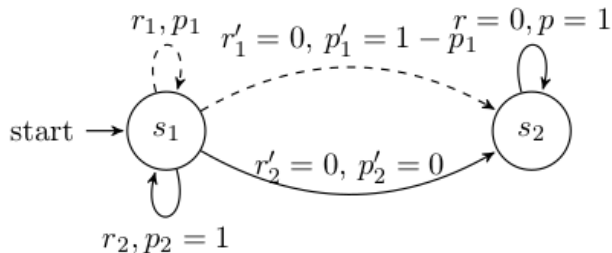
- $Q(s_1, a_i) = \frac{r_i}{1-\gamma p_i}$ ,  $i = 1, 2$ .
- Can easily construct  $\psi$  and  $\bar{\psi}$  such that:
  - Both  $\psi$  and  $\bar{\psi}$  satisfy  $\frac{r_1}{1-\gamma p_1} > \frac{r_2}{1-\gamma p_2}$ .

# IT Lower bound: Intractable!



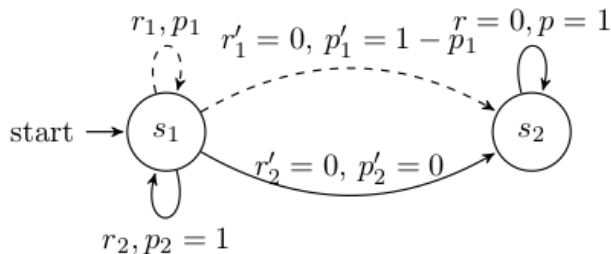
- $Q(s_1, a_i) = \frac{r_i}{1-\gamma p_i}$ ,  $i = 1, 2$ .
- Can easily construct  $\psi$  and  $\bar{\psi}$  such that:
  - Both  $\psi$  and  $\bar{\psi}$  satisfy  $\frac{r_1}{1-\gamma p_1} > \frac{r_2}{1-\gamma p_2}$ .
  - $\phi = \frac{\psi + \bar{\psi}}{2}$  satisfies  $\frac{r_1}{1-\gamma p_1} < \frac{r_2}{1-\gamma p_2}$ .

# IT Lower bound: Intractable!



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

# IT Lower bound: Intractable!



- $\text{Alt}(\phi)$  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

	<b>MAB</b>	<b>MDP</b>
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
IT lower bound	Tractable	Not Tractable

# Upper bound: Idea

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

# Upper bound: Idea

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

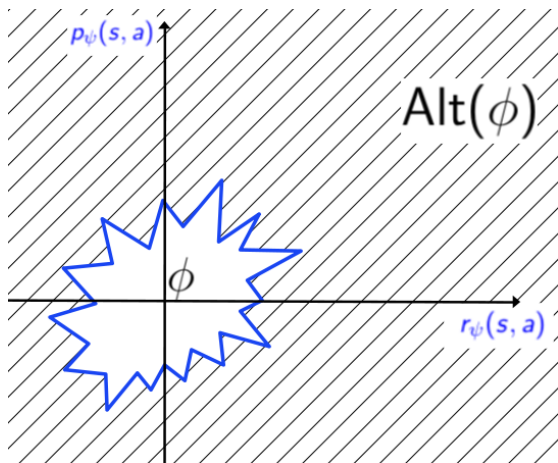


Figure:  $\text{Alt}(\phi)$ : Non-convex boundary

# Upper bound: Idea

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

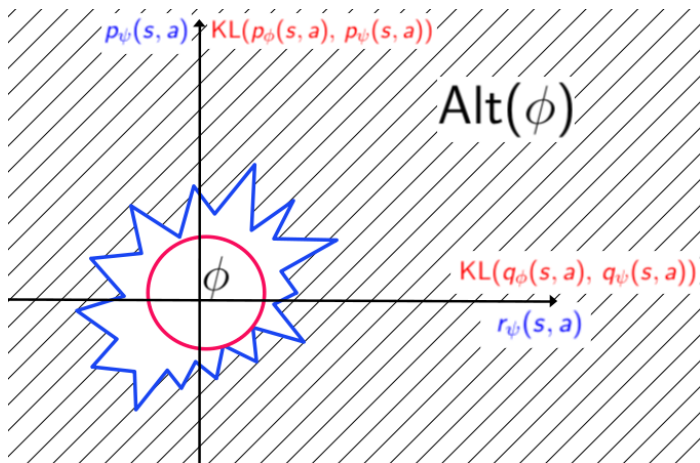


Figure: KL Ball

# Upper bound of the characteristic time

## Theorem 1 (Upper bound of minimal sample complexity)

For all vectors  $\omega$  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

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

# Upper bound of $T(\phi, \omega)$ : sketch of the proof

- Using Lemma 2:

$$T(\phi, \omega)^{-1} = \min_{s, a \neq \pi^*(s)} \inf_{\psi \in \text{Alt}_{sa}(\phi)} \omega_{sa} \text{KL}_{\phi|\psi}(s, a) + \sum_x \omega_{x, \pi^*(x)} \text{KL}_{\phi|\psi}(x, \pi^*(x)).$$

where  $\text{Alt}_{sa}(\phi) \triangleq \{\psi : Q_{\psi}^{\pi^*}(s, a) > V_{\psi}^{\pi^*}(s)\}$ .

# Upper bound of $T(\phi, \omega)$ : sketch of the proof

- Using Lemma 2:

$$T(\phi, \omega)^{-1} = \min_{s, a \neq \pi^*(s)} \inf_{\psi \in \text{Alt}_{sa}(\phi)} \omega_{sa} \text{KL}_{\phi|\psi}(s, a) + \sum_x \omega_{x, \pi^*(x)} \text{KL}_{\phi|\psi}(x, \pi^*(x)).$$

where  $\text{Alt}_{sa}(\phi) \triangleq \{\psi : Q_{\psi}^{\pi^*}(s, a) > V_{\psi}^{\pi^*}(s)\}$ .

- Introduce the suboptimality gaps:  $\Delta_{sa} \triangleq V_{\phi}^{\pi^*}(s) - Q_{\phi}^{\pi^*}(s, a)$ :

$$(Q_{\psi}^{\pi^*} - Q_{\phi}^{\pi^*})(s, a) - (V_{\psi}^{\pi^*} - V_{\phi}^{\pi^*})(s) > \Delta_{sa}.$$

# Upper bound of $T(\phi, \omega)$ : sketch of the proof

- Using Lemma 2:

$$T(\phi, \omega)^{-1} = \min_{s, a \neq \pi^*(s)} \inf_{\psi \in \text{Alt}_{sa}(\phi)} \omega_{sa} \text{KL}_{\phi|\psi}(s, a) + \sum_x \omega_{x, \pi^*(x)} \text{KL}_{\phi|\psi}(x, \pi^*(x)).$$

where  $\text{Alt}_{sa}(\phi) \triangleq \{\psi : Q_{\psi}^{\pi^*}(s, a) > V_{\psi}^{\pi^*}(s)\}$ .

- Introduce the suboptimality gaps:  $\Delta_{sa} \triangleq V_{\phi}^{\pi^*}(s) - Q_{\phi}^{\pi^*}(s, a)$ :

$$(Q_{\psi}^{\pi^*} - Q_{\phi}^{\pi^*})(s, a) - (V_{\psi}^{\pi^*} - V_{\phi}^{\pi^*})(s) > \Delta_{sa}.$$

- Rewrite the condition in terms of the differences in kernels:

$$dr(s, a) + \gamma dp(s, a)^T V_{\phi}^{\pi^*} + \gamma [p_{\psi}(s, a) - \mathbb{1}(s)] dV^{\pi^*} > \Delta_{sa}$$

where

$$dr(s, a) = r_{\psi}(s, a) - r_{\phi}(s, a), dV^{\pi^*}(s, a) = V_{\psi}^{\pi^*}(s, a) - V_{\phi}^{\pi^*}(s, a) \text{ etc}$$



## Upper bound of $T(\phi, \omega)$ : sketch of the proof

- $dr(s, a) + \gamma dp(s, a)^T V_{\phi}^{\pi^*} + \gamma [p_{\psi}(s, a) - \mathbb{1}(s)] dV^{\pi^*} > \Delta_{sa}$

## Upper bound of $T(\phi, \omega)$ : sketch of the proof

- $dr(s, a) + \gamma dp(s, a)^T V_\phi^{\pi^*} + \gamma[p_\psi(s, a) - \mathbb{1}(s)]dV^{\pi^*} > \Delta_{sa}$
- $dV^{\pi^*} = A + B$  where:

$$A \triangleq \left( I - \gamma P_\psi^{\pi^*} \right)^{-1} \left[ r_\psi^{\pi^*} - r_\phi^{\pi^*} \right].$$

$$B \triangleq \left[ \left( I - \gamma P_\psi^{\pi^*} \right)^{-1} - \left( I - \gamma P_\phi^{\pi^*} \right)^{-1} \right] r_\phi^{\pi^*}.$$

## Upper bound of $T(\phi, \omega)$ : sketch of the proof

- $dr(s, a) + \gamma dp(s, a)^T V_\phi^{\pi^*} + \gamma[p_\psi(s, a) - \mathbb{1}(s)]dV^{\pi^*} > \Delta_{sa}$
- $dV^{\pi^*} = A + B$  where:

$$A \triangleq \left( I - \gamma P_\psi^{\pi^*} \right)^{-1} \left[ r_\psi^{\pi^*} - r_\phi^{\pi^*} \right].$$

$$B \triangleq \left[ \left( I - \gamma P_\psi^{\pi^*} \right)^{-1} - \left( I - \gamma P_\phi^{\pi^*} \right)^{-1} \right] r_\phi^{\pi^*}.$$

- $dr(s, a) = \alpha_1 \Delta_{sa}$ ,  $dp(s, a) = \alpha_2 \Delta_{sa}$ ,  $A = \alpha_3 \Delta_{sa}$ ,  $B = \alpha_4 \Delta_{sa}$ , with  $\sum_i \alpha_i > 1$

# Upper bound of $T(\phi, \omega)$ : sketch of the proof

- $dr(s, a) + \gamma dp(s, a)^T V_\phi^{\pi^*} + \gamma[p_\psi(s, a) - \mathbb{1}(s)]dV^{\pi^*} > \Delta_{sa}$
- $dV^{\pi^*} = A + B$  where:

$$A \triangleq \left(I - \gamma P_\psi^{\pi^*}\right)^{-1} \left[r_\psi^{\pi^*} - r_\phi^{\pi^*}\right].$$

$$B \triangleq \left[\left(I - \gamma P_\psi^{\pi^*}\right)^{-1} - \left(I - \gamma P_\phi^{\pi^*}\right)^{-1}\right] r_\phi^{\pi^*}.$$

- $dr(s, a) = \alpha_1 \Delta_{sa}$ ,  $dp(s, a) = \alpha_2 \Delta_{sa}$ ,  $A = \alpha_3 \Delta_{sa}$ ,  $B = \alpha_4 \Delta_{sa}$ , with  $\sum_i \alpha_i > 1$
- Use Pinsker inequality and transportation lemmas to relate  $dr$ ,  $dp$  to  $KL(q_\phi, q_\psi)$ ,  $KL(p_\phi, p_\psi)$ :

$$\frac{1}{2}(\alpha_1 \Delta_{sa})^2 \leq KL(q_\phi(s, a), q_\psi(s, a)).$$

# Upper bound: sketch of the proof

- Use IT inequalities to relate  $dr, dp$  to  $KL(q_\phi, q_\psi), KL(p_\phi, p_\psi)$ :

$$\frac{\alpha_1^2}{T_1} \leq KL(q_\phi(\cdot|s, a), q_\psi(s, a)).$$

$$\frac{\alpha_2^2}{T_2} \leq KL(p_\phi(\cdot|s, a), p_\psi(s, a)).$$

$$\frac{\alpha_3^2}{T_3} \leq \max_s KL(q_\phi(\cdot|s, \pi^*(s)), q_\psi(\cdot|s, \pi^*(s))).$$

$$\frac{\alpha_4^2}{T_4} \leq \max_s KL(p_\phi(s, \pi^*(s)), p_\psi(s, \pi^*(s))).$$

## Upper bound: sketch of the proof

- Use IT inequalities to relate  $dr, dp$  to  $KL(q_\phi, q_\psi), KL(p_\phi, p_\psi)$ :

$$\frac{\alpha_1^2}{T_1} \leq KL(q_\phi(\cdot|s, a), q_\psi(s, a)).$$

$$\frac{\alpha_2^2}{T_2} \leq KL(p_\phi(\cdot|s, a), p_\psi(s, a)).$$

$$\frac{\alpha_3^2}{T_3} \leq \max_s KL(q_\phi(\cdot|s, \pi^*(s)), q_\psi(\cdot|s, \pi^*(s))).$$

$$\frac{\alpha_4^2}{T_4} \leq \max_s KL(p_\phi(s, \pi^*(s)), p_\psi(s, \pi^*(s))).$$

- Sum-up the bounds and optimize over  $\alpha$ .

# Upper bound: sketch of the proof

- Use IT inequalities to relate  $dr, dp$  to  $KL(q_\phi, q_\psi), KL(p_\phi, p_\psi)$ :

$$\frac{\alpha_1^2}{T_1} \leq KL(q_\phi(\cdot|s, a), q_\psi(s, a)).$$

$$\frac{\alpha_2^2}{T_2} \leq KL(p_\phi(\cdot|s, a), p_\psi(s, a)).$$

$$\frac{\alpha_3^2}{T_3} \leq \max_s KL(q_\phi(\cdot|s, \pi^*(s)), q_\psi(\cdot|s, \pi^*(s))).$$

$$\frac{\alpha_4^2}{T_4} \leq \max_s KL(p_\phi(s, \pi^*(s)), p_\psi(s, \pi^*(s))).$$

- Sum-up the bounds and optimize over  $\alpha$ .
- Gives a lower bound of  $T(\phi, \omega)^{-1} =$

$$\min_{s, a \neq \pi^*(s)} \inf_{\psi \in \text{Alt}_{sa}(\phi)} \omega_{sa} KL_{\phi|\psi}(s, a) + \sum_x \omega_{x, \pi^*(x)} KL_{\phi|\psi}(x, \pi^*(x)).$$

# KLB-TS: Sampling rule



- 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 !

- $\bar{\omega}_{sa} \propto \frac{1 + \text{Var}_{p_{\phi}(s,a)}[V_{\phi}^*]}{\Delta_{s,a}^2}$ .

- 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 !

- $\bar{\omega}_{sa} \propto \frac{1 + \text{Var}_{p_{\phi}(s,a)}[V_{\phi}^*]}{\Delta_{s,a}^2}$ .
- $\bar{\omega}_{s, \pi^*(s)} \propto \frac{1 + \text{Var}_{\max}^*[V_{\phi}^*]}{\Delta_{\min}^2 (1-\gamma)^2}$ .

- 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 !**

- Use C-Tracking [Garivier and Kaufmann, 2016]:

- 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 !**

- Use C-Tracking [Garivier and Kaufmann, 2016]:
  - Project  $\bar{\omega}(\hat{\phi}_t)$  on  $\{\omega \in \Sigma : \forall (s, a), \omega_{sa} \geq \frac{1}{\sqrt{t}}\}$  to get  $\tilde{\omega}(\hat{\phi}_t)$ .

- 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** !

- Use C-Tracking [Garivier and Kaufmann, 2016]:
  - Project  $\bar{\omega}(\hat{\phi}_t)$  on  $\{\omega \in \Sigma : \forall (s, a), \omega_{sa} \geq \frac{1}{\sqrt{t}}\}$  to get  $\tilde{\omega}(\hat{\phi}_t)$ .
  - $(s_{t+1}, a_{t+1}) \in \arg \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sum_{s=1}^t \tilde{\omega}_{sa}(\hat{\phi}_s) - N_{sa}(t)$ .

- 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** !

- Use C-Tracking [Garivier and Kaufmann, 2016]:
  - Project  $\bar{\omega}(\hat{\phi}_t)$  on  $\{\omega \in \Sigma : \forall (s, a), \omega_{sa} \geq \frac{1}{\sqrt{t}}\}$  to get  $\tilde{\omega}(\hat{\phi}_t)$ .
  - $(s_{t+1}, a_{t+1}) \in \arg \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sum_{s=1}^t \tilde{\omega}_{sa}(\hat{\phi}_s) - N_{sa}(t)$ .
- 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

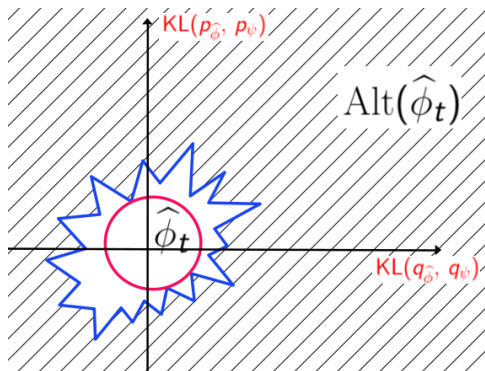


Figure: KL-Ball Stopping rule



# KLB-TS: stopping rule

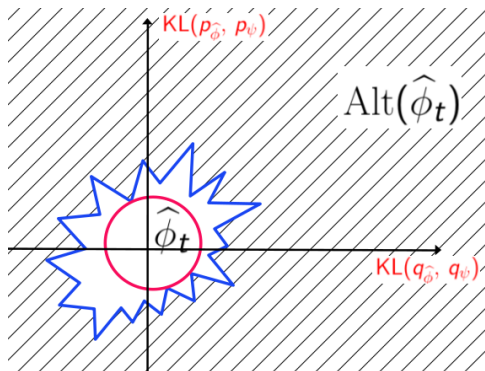


Figure: KL-Ball Stopping rule

- Need to ensure that  $\phi$  falls within the KL-ball with probability  $1 - \delta$ .

## Theorem 3

KLB-TS has a sample complexity  $\tau_\delta$  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:

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

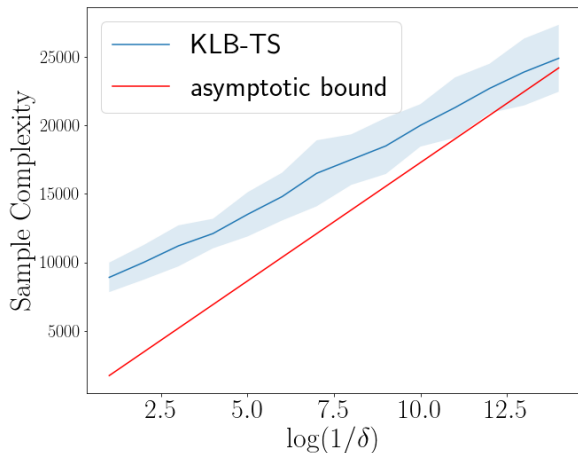


Figure: Asymptotic bound:  $S=A=2$ ,  $\gamma = 0.5$ .

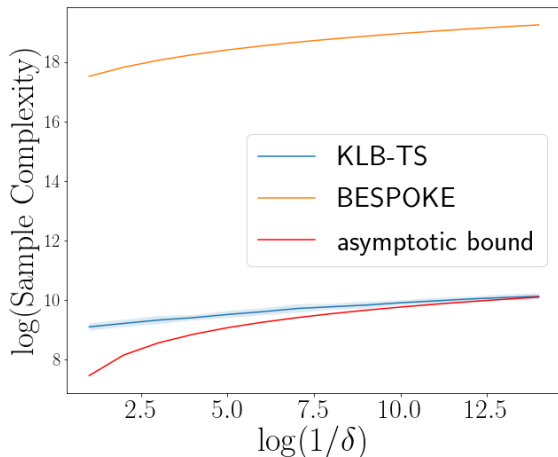


Figure: KLB-TS vs. BESPOKE.  $S=A=2$ ,  $\gamma = 0.5$ .

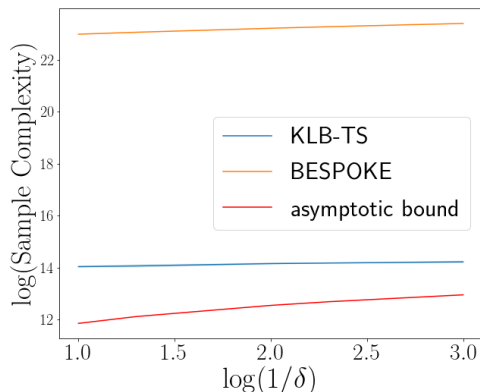






Figure: KLB-TS vs. BESPOKE.  $S = 5, A = 10, \gamma = 0.7$ .




- Most of BESPOKE's sample complexity comes from the burn-in phase  $\Omega\left(\frac{S^2 A \log(1/\delta)}{(1-\gamma)^2}\right)$ .

- 1 Algorithms designed using *problem-specific* bounds can achieve better sample complexity than minimax ones.
- 2 Contrary to MAB, IT lower bound is intractable for MDPs.
- 3 We can derive problem-specific surrogates which :
  - Are *explicit*, depending on functionals of the MDP.
  - Have a corresponding allocation that is easy to compute.
- 4 Can be used to devise (Asymptotically) Matching algorithm.
- 5 First step towards understanding problem-specific  $\varepsilon$ -optimal policy identification.




Thanks !

-  Agarwal, A., Kakade, S., and Yang, L. F. (2020). Model-based reinforcement learning with a generative model is minimax optimal. volume 125 of *Proceedings of Machine Learning Research*, pages 67–83. PMLR.
-  Azar, M. G., Munos, R., and Kappen, H. J. (2013). Minimax pac bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):325–349.
-  Even-Dar, E. and Mansour, Y. (2003). Learning rates for q-learning. *Journal of Machine Learning Research*, 5(Dec):1–25.
-  Gabillon, V., Ghavamzadeh, M., and Lazaric, A. (2012). Best arm identification: A unified approach to fixed budget and fixed confidence. In *NIPS*.



-  Garivier, A. and Kaufmann, E. (2016). Optimal best arm identification with fixed confidence. In Feldman, V., Rakhlin, A., and Shamir, O., editors, *29th Annual Conference on Learning Theory*, volume 49 of *Proceedings of Machine Learning Research*, pages 998–1027, Columbia University, New York, New York, USA. PMLR.
-  Kalyanakrishnan, S., Tewari, A., Auer, P., and Stone, P. (2012). Pac subset selection in stochastic multi-armed bandits. In *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, pages 655–662, New York, NY, USA. ACM.
-  Kearns, M. and Singh, S. (1999). Finite-sample convergence rates for q-learning and indirect algorithms. *Advances in Neural Information Processing*, 11.

## References III

-  Li, G., Wei, Y., Chi, Y., Gu, Y., and Chen, Y. (2020). Breaking the sample size barrier in model-based reinforcement learning with a generative model. *arXiv preprint, arXiv:2005.12900*.
-  Sidford, A., Wang, M., Wu, X., Yang, L., and Ye, Y. (2018). Near-optimal time and sample complexities for solving markov decision processes with a generative model. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems 31*, pages 5186–5196. Curran Associates, Inc.
-  Zanette, A., Kochenderfer, M. J., and Brunskill, E. (2019). Almost horizon-free structure-aware best policy identification with a generative model. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché Buc, F., Fox, E., and Garnett, R., editors, *Advances in Neural Information Processing Systems 32*, pages 5625–5634. Curran Associates, Inc.