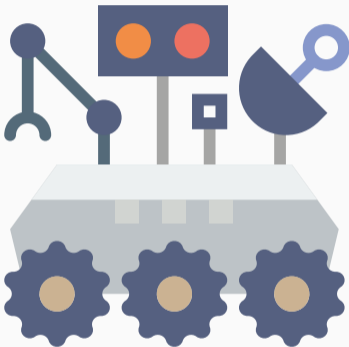# Optimal Policy Identification

Wouter M. Koolen

CWI
Centrum Wiskunde & Informatica

UNIVERSITY
OF TWENTE.

RL Seminar
University of Twente
Friday 9th June, 2023

A. Al Marjani and A. Proutiere (2021). "Adaptive sampling for best policy identification in Markov decision processes". In: **International Conference on Machine Learning**. PMLR, pp. 7459–7468

## Outline

1. Setup

# Markov Decision Process

Finite sets of states $S$ and **actions** $A$. Rewards bounded in $[0, 1]$.

## Markov Decision Process

Finite sets of states $S$ and **actions** $A$. Rewards bounded in $[0, 1]$.

MDP $\phi = (p_\phi, q_\phi)$ specified by

- dynamics $p_\phi(s'|s, a)$ and
- rewards $q_\phi(r|s, a)$.

Finite sets of states $S$ and **actions** $A$. Rewards bounded in $[0, 1]$.

MDP $\phi = (p_\phi, q_\phi)$ specified by

- dynamics $p_\phi(s'|s, a)$ and
- rewards $q_\phi(r|s, a)$.

We will write $r_\phi(s, a)$ for the mean reward (of doing $a$ in $s$ under $\phi$).

A **policy** is map $\pi : S \to A$.

Executing a policy $\pi$ from state $s$ under $\phi$ gives a sequence $(s_t^\pi)_{t \geq 0}$ with

$$s_0^\pi \;=\; s \qquad \text{and} \qquad s_{t+1}^\pi \;\sim\; p_\phi\big(\cdot \big| s_t^\pi, \pi(s_t^\pi)\big)$$

A **policy** is map $\pi : S \to A$.

Executing a policy $\pi$ from state $s$ under $\phi$ gives a sequence $(s_t^\pi)_{t \geq 0}$ with

$$s_0^\pi = s \qquad \text{and} \qquad s_{t+1}^\pi \sim p_\phi\left(\cdot \big| s_t^\pi, \pi(s_t^\pi)\right)$$

Let's introduce discount factor $\gamma \in (0, 1)$.

Value function of policy $\pi$ in $\phi$:

$$V_\phi^\pi(s) = \mathbb{E}_\phi\left[\sum_{t=0}^{\infty} \gamma^t r_\phi\left(s_t^\pi, \pi(s_t^\pi)\right) \bigg| s_0^\pi = s\right]$$

A **policy** is map $\pi : S \to A$.

Executing a policy $\pi$ from state $s$ under $\phi$ gives a sequence $(s_t^\pi)_{t \geq 0}$ with

$$s_0^\pi = s \qquad \text{and} \qquad s_{t+1}^\pi \sim p_\phi\big(\cdot \big| s_t^\pi, \pi(s_t^\pi)\big)$$

Let's introduce discount factor $\gamma \in (0, 1)$.

Value function of policy $\pi$ in $\phi$:

$$V_\phi^\pi(s) = \mathbb{E}_\phi \left[ \sum_{t=0}^{\infty} \gamma^t r_\phi\left(s_t^\pi, \pi(s_t^\pi)\right) \middle| s_0^\pi = s \right]$$

**Optimal policy** for $\phi$ is $\pi^*(\phi) \coloneqq \underset{\pi : S \to A}{\arg\max} \, V_\phi^\pi$ \quad (NB: not a scalar(!))

**Problem**

*Given any unknown MDP $\phi$, identify its optimal policy $\pi^*(\phi)$ from interactive exploration.*

We want

- Reliability: output policy is indeed optimal
- **Efficiency**: as few samples as possible

Fix unknown MDP $\phi$.

**Protocol (generative model)**

**for** $t = 1, 2, \ldots, \tau$

- Learner picks state $s_t$ and action $a_t$
- Learner observes reward $r_t \sim q_\phi(\cdot | s_t, a_t)$ and successor state $s'_t \sim p_\phi(\cdot | s_t, a_t)$

Learner recommends policy $\hat{\pi}$.

Fix unknown MDP $\phi$.

**Protocol (generative model)**

**for** $t = 1, 2, \ldots, \tau$
- Learner picks state $s_t$ and action $a_t$
- Learner observes reward $r_t \sim q_\phi(\cdot | s_t, a_t)$ and successor state $s'_t \sim p_\phi(\cdot | s_t, a_t)$

Learner recommends policy $\hat{\hat{\pi}}$.

NB: generative model different from **navigation**, where $s_{t+1} = s'_t$.

Fix unknown MDP $\phi$.

### Protocol (generative model)

**for** $t = 1, 2, \ldots, \tau$
- Learner picks state $s_t$ and action $a_t$
- Learner observes reward $r_t \sim q_\phi(\cdot | s_t, a_t)$ and successor state $s'_t \sim p_\phi(\cdot | s_t, a_t)$

Learner recommends policy $\hat{\pi}$.

NB: generative model different from **navigation**, where $s_{t+1} = s'_t$.

### Definition (Correctness)

Fix confidence $\delta \in (0, 1)$. A learner is $\delta$-correct if $\mathbb{P}_\phi(\hat{\pi} \neq \pi^*(\phi)) \leq \delta$.

## Interactive Exploration: Interaction with a generative model

Fix unknown MDP $\phi$.

**Protocol (generative model)**

**for** $t = 1, 2, \ldots, \tau$
- Learner picks state $s_t$ and action $a_t$
- Learner observes reward $r_t \sim q_\phi(\cdot | s_t, a_t)$ and successor state $s'_t \sim p_\phi(\cdot | s_t, a_t)$

Learner recommends policy $\hat{\pi}$.

NB: generative model different from **navigation**, where $s_{t+1} = s'_t$.

**Definition (Correctness)**

Fix confidence $\delta \in (0, 1)$. A learner is $\delta$-correct if $\mathbb{P}_\phi(\hat{\pi} \neq \pi^*(\phi)) \leq \delta$.

**Quest: efficiency**

Minimise sample complexity $\phi \mapsto \mathbb{E}_\phi[\tau]$ over all $\delta$-correct learners.

Uniform sampling: sample each pair $(s, a)$ for $n$ times. Concentration results say

- Estimate $r_\phi(s, a)$ up to precision $1/\sqrt{n}$.
- Estimate $p_\phi(s'|s, a)$ up to precision $1/\sqrt{n}$ for each $s'$.

Uniform sampling: sample each pair $(s, a)$ for $n$ times. Concentration results say

- Estimate $r_\phi(s, a)$ up to precision $1/\sqrt{n}$.
- Estimate $p_\phi(s'|s, a)$ up to precision $1/\sqrt{n}$ for each $s'$.

The value function of $\pi$ satisfies the recurrence

$$V_\phi^\pi(s) = r_\phi(s, \pi(s)) + \gamma \sum_{s'} p_\phi(s'|s, \pi(s)) V_\phi^\pi(s')$$

That is

$$V_\phi^\pi = \left(I - \gamma \sum_{s,s'} p_\phi(s'|s, \pi(s)) e_s e_{s'}^\mathsf{T}\right)^{-1} \sum_s e_s r_\phi(s, \pi(s))$$

A perturbation argument gives $\|V_\phi^\pi - V_{\hat\phi}^\pi\|_\infty \leq O\left(\frac{1}{\sqrt{n}(1-\gamma)^2}\right)$.

Uniform sampling: sample each pair $(s, a)$ for $n$ times. Concentration results say

- Estimate $r_\phi(s, a)$ up to precision $1/\sqrt{n}$.
- Estimate $p_\phi(s'|s, a)$ up to precision $1/\sqrt{n}$ for each $s'$.

The value function of $\pi$ satisfies the recurrence

$$V_\phi^\pi(s) \;=\; r_\phi(s, \pi(s)) + \gamma \sum_{s'} p_\phi(s'|s, \pi(s)) \, V_\phi^\pi(s')$$

That is

$$V_\phi^\pi \;=\; \left( I - \gamma \sum_{s,s'} p_\phi(s'|s, \pi(s)) e_s e_{s'}^\mathsf{T} \right)^{-1} \sum_s e_s r_\phi(s, \pi(s))$$

A perturbation argument gives $\|V_\phi^\pi - V_{\hat\phi}^\pi\|_\infty \le O\left( \frac{1}{\sqrt{n}(1-\gamma)^2} \right)$.

Obtain $\epsilon$-optimal policy in using $O\left( \frac{SA}{\epsilon^2(1-\gamma)^4} \ln(SA) \right)$ samples.

**Theorem (Azar, Munos, and Kappen, 2013)**

*Can identify an $\epsilon$-optimal policy in samples*

$$O\left(\frac{SA}{\epsilon^2(1-\gamma)^3}\ln(SA)\right)$$

And no algorithm can do uniformly better. Worst-case optimal!

**Theorem (Azar, Munos, and Kappen, 2013)**

*Can identify an $\epsilon$-optimal policy in samples*

$$O\left(\frac{SA}{\epsilon^2(1-\gamma)^3}\ln(SA)\right)$$

And no algorithm can do uniformly better. Worst-case optimal!

But maybe some MDPs $\phi$ are **easy**?

**Theorem (Azar, Munos, and Kappen, 2013)**

*Can identify an $\epsilon$-optimal policy in samples*

$$O\left(\frac{SA}{\epsilon^2(1-\gamma)^3}\ln(SA)\right)$$

And no algorithm can do uniformly better. Worst-case optimal!

But maybe some MDPs $\phi$ are **easy**?

Can we have instance-dependent / instance-optimal results?

## Discrimination

In MDP $\phi$, the average total evidence collected against MDP $\psi$ is

$$
\begin{aligned}
&\mathsf{KL}_{\phi|\psi}\big((s_t, a_t, s_t', r_t)_{t=1}^{\tau}, \hat{\pi}\big) \\
&= \sum_{s,a} \mathbb{E}_\phi[N_{s,a}(\tau)]\Big\{\mathsf{KL}\big(p_\phi(\cdot|s,a)\big\|p_\psi(\cdot|s,a)\big) + \mathsf{KL}\big(q_\phi(\cdot|s,a)\big\|q_\psi(\cdot|s,a)\big)\Big\} \\
&= \mathbb{E}_\phi[\tau] \sum_{s,a} \frac{\mathbb{E}_\phi[N_{s,a}(\tau)]}{\mathbb{E}_\phi[\tau]}\Big\{\mathsf{KL}\big(p_\phi(\cdot|s,a)\big\|p_\psi(\cdot|s,a)\big) + \mathsf{KL}\big(q_\phi(\cdot|s,a)\big\|q_\psi(\cdot|s,a)\big)\Big\}
\end{aligned}
$$

## Discrimination

In MDP $\phi$, the average total evidence collected against MDP $\psi$ is

$$
\begin{aligned}
& \mathsf{KL}_{\phi|\psi}\big((s_t, a_t, s_t', r_t)_{t=1}^{\tau}, \hat{\pi}\big) \\
&= \sum_{s,a} \mathbb{E}_\phi[N_{s,a}(\tau)]\Big\{\mathsf{KL}\big(p_\phi(\cdot|s,a)\big\|p_\psi(\cdot|s,a)\big) + \mathsf{KL}\big(q_\phi(\cdot|s,a)\big\|q_\psi(\cdot|s,a)\big)\Big\} \\
&= \mathbb{E}_\phi[\tau] \sum_{s,a} \frac{\mathbb{E}_\phi[N_{s,a}(\tau)]}{\mathbb{E}_\phi[\tau]}\Big\{\mathsf{KL}\big(p_\phi(\cdot|s,a)\big\|p_\psi(\cdot|s,a)\big) + \mathsf{KL}\big(q_\phi(\cdot|s,a)\big\|q_\psi(\cdot|s,a)\big)\Big\}
\end{aligned}
$$

Moreover, if Learner is $\delta$-correct and $\pi^*(\phi) \neq \pi^*(\psi)$ then

$$
\mathsf{KL}_{\phi|\psi}\big((s_t, a_t, s_t', r_t)_{t=1}^{\tau}, \hat{\pi}\big) \;\geq\; \mathsf{KL}_{\phi|\psi}\big(\mathbf{1}_{\hat{\pi}=\pi^*(\phi)}\big) \;\geq\; \mathsf{KL}(1-\delta, \delta) \;\approx\; \ln\frac{1}{\delta}
$$

## Discrimination

In MDP $\phi$, the average total evidence collected against MDP $\psi$ is

$$
\begin{aligned}
& \mathsf{KL}_{\phi|\psi}\big((s_t, a_t, s'_t, r_t)_{t=1}^\tau, \hat\pi\big) \\
&= \sum_{s,a} \mathbb{E}_\phi[N_{s,a}(\tau)]\Big\{\mathsf{KL}\big(p_\phi(\cdot|s,a)\big\|p_\psi(\cdot|s,a)\big) + \mathsf{KL}\big(q_\phi(\cdot|s,a)\big\|q_\psi(\cdot|s,a)\big)\Big\} \\
&= \mathbb{E}_\phi[\tau]\sum_{s,a} \frac{\mathbb{E}_\phi[N_{s,a}(\tau)]}{\mathbb{E}_\phi[\tau]}\Big\{\mathsf{KL}\big(p_\phi(\cdot|s,a)\big\|p_\psi(\cdot|s,a)\big) + \mathsf{KL}\big(q_\phi(\cdot|s,a)\big\|q_\psi(\cdot|s,a)\big)\Big\}
\end{aligned}
$$

Moreover, if Learner is $\delta$-correct and $\pi^*(\phi) \neq \pi^*(\psi)$ then

$$
\mathsf{KL}_{\phi|\psi}\big((s_t, a_t, s'_t, r_t)_{t=1}^\tau, \hat\pi\big) \;\geq\; \mathsf{KL}_{\phi|\psi}\big(\mathbf{1}_{\hat\pi=\pi^*(\phi)}\big) \;\geq\; \mathsf{KL}(1-\delta, \delta) \;\approx\; \ln\frac{1}{\delta}
$$

So all in all

$$
\mathbb{E}_\phi[\tau] \;\geq\; \frac{\ln\frac{1}{\delta}}{\displaystyle\max_{\boldsymbol{w}\in\triangle_{SA}} \min_{\psi:\pi^*(\psi)\neq\pi^*(\phi)} \sum_{s,a} w_{s,a}\Big\{\mathsf{KL}\big(p_\phi(\cdot|s,a)\big\|p_\psi(\cdot|s,a)\big) + \mathsf{KL}\big(q_\phi(\cdot|s,a)\big\|q_\psi(\cdot|s,a)\big)\Big\}}
$$

## Optimal Algorithm (Track-and-Stop template)

Can we have a single algorithm so that for all $\phi$,

$$\mathbb{E}_\phi[\tau] \leq \frac{\ln\frac{1}{\delta}}{\max\limits_{\boldsymbol{w}\in\triangle_{SA}} \min\limits_{\psi:\pi^*(\psi)\neq\pi^*(\phi)} \sum\limits_{s,a} w_{s,a}\Big\{\mathsf{KL}\big(p_\phi(\cdot|s,a)\big\|p_\psi(\cdot|s,a)\big) + \mathsf{KL}\big(q_\phi(\cdot|s,a)\big\|q_\psi(\cdot|s,a)\big)\Big\}}+\text{tiny?}$$

## Optimal Algorithm (Track-and-Stop template)

Can we have a single algorithm so that for all $\phi$,

$$\mathbb{E}_\phi[\tau] \leq \frac{\ln\frac{1}{\delta}}{\max\limits_{\boldsymbol{w}\in\triangle_{SA}} \min\limits_{\psi:\pi^*(\psi)\neq\pi^*(\phi)} \sum\limits_{s,a} w_{s,a}\Big\{\mathsf{KL}\big(p_\phi(\cdot|s,a)\big\|p_\psi(\cdot|s,a)\big) + \mathsf{KL}\big(q_\phi(\cdot|s,a)\big\|q_\psi(\cdot|s,a)\big)\Big\}} + \text{tiny?}$$

### Track and Stop

For $t = 1, 2, \ldots$
- Form estimate $\hat{\phi}_t$ of MDP
- Compute oracle weights $\boldsymbol{w}^*(\hat{\phi}_t)$.
- Track $\boldsymbol{w}^*$ with sub-linear forced exploration.
- Stop/recommend using GLRT (generalised likelihood ratio test).

Analysis: as $t \to \infty$, then $\hat{\phi}_t \to \phi$ hence $\boldsymbol{w}^*(\hat{\phi}_t) \to \boldsymbol{w}^*(\phi)$ and we stop at optimal time + tiny.

The instance dependent problem complexity of MDPs is (apparently)

$$\frac{1}{\max\limits_{\boldsymbol{w}\in\triangle_{SA}} \min\limits_{\psi:\pi^*(\psi)\neq\pi^*(\phi)} \sum\limits_{s,a} w_{s,a}\Big\{\mathsf{KL}\big(p_\phi(\cdot|s,a)\big\|p_\psi(\cdot|s,a)\big) + \mathsf{KL}\big(q_\phi(\cdot|s,a)\big\|q_\psi(\cdot|s,a)\big)\Big\}}$$

as lower and upper bounds match.

Algorithmics **not settled**.

Current algorithms target relaxations instead.

Thanks!

📄 Al Marjani, A. and A. Proutiere (2021). "Adaptive sampling for best policy identification in Markov decision processes". In: **International Conference on Machine Learning**. PMLR, pp. 7459–7468.

📄 Azar, M. G., R. Munos, and H. J. Kappen (2013). "Minimax PAC bounds on the sample complexity of reinforcement learning with a generative model". In: **Machine learning** 91, pp. 325–349.