

Ergodic Blind MDPs: Decidability of Approximation



K. Chatterjee¹



R. Saona¹



B. Ziliotto^{2,3}



D. Lurie²

¹ISTA

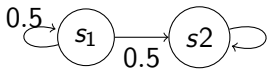
²CEREMADE, Université Paris Dauphine,

³CNRS, PSL Research Institute

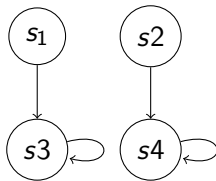
For ergodic blind MDPs with \liminf average objective,

- **COMPUTING** the value is **UNDECIDABLE**.
- **APPROXIMATING** the value is **DECIDABLE**.

Blind MDPs

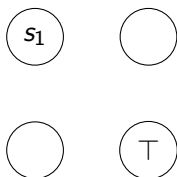


Action *wait*



Action *go*

Reachability objective



We focus in the maximum probability of reaching the target

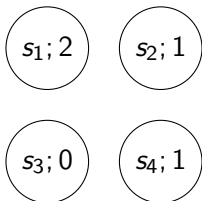
$$\sup_{\sigma} \mathbb{P}^{\sigma}(\exists n \ S_n = T)$$

Theorem (Madani et. al. 2003)

Approximating *the value of a blind MDP with reachability objective is* **undecidable**.

Via a reduction of the acceptance of the empty word in a two-counter Turing Machine.

Liminf average objective



We focus in the maximum expected longrun average reward

$$\sup_{\sigma} \mathbb{E}^{\sigma} \left(\liminf_{n \rightarrow \infty} \frac{1}{n} \sum_{m \geq n} r(S_m) \right)$$

Corollary (Madani et. al. 2003)

Approximating *the value of a blind MDP with liminf average objective is **undecidable**.*

Theorem (Saona et. al. 2021)

*Consider blind MDPs with liminf average objective. For all $\varepsilon > 0$, there exists a **finite-recall ε -optimal strategy**.*

Via an ergodic theorem and block strategies.

By the **undeciability** of approximation, there is **no explicit bound** on the amount of recall needed to achieve ε -optimality.

What **subclasses** of blind MDPs **allow a bound** on the recall?

An explicit bound leads to decidability of approximation.

An **ergodic** blind MDP approximately **forgets the initial belief** after long time.

Formally, for all $\varepsilon > 0$, there exists $n_0 \geq 1$ such that, for all $n \geq n_0$,

$$\mathbb{P}_{p_1}^\sigma(S_n = \cdot) \approx_\varepsilon \mathbb{P}_{q_1}^\sigma(S_n = \cdot)$$

for all σ, p_1, q_1 .

Lemma

Checking if a given blind MDP is ergodic is in EXPSPACE.

Theorem

For ergodic blind MDPs with liminf average objective, approximating the value is decidable.

- Consider $\varepsilon > 0$.
- After n the belief is ε -“independent” of the initial belief.
- After n step approximate the belief.
- Since approximations are “independent”, the error does not accumulate over time.

Theorem

*For ergodic blind MDPs with \liminf average objective, **computing the value is undecidable.***

- Consider a PFA.
- Construct an ergodic blind MDP by leaking to a state and allowing restart.
 - $\text{value}(\text{blind MDP}) > 1/2 \iff \text{value}(\text{PFA}) > 1/2$

Proposition

All ergodic blind MDPs with reachability objective have value one.

Starting from the target state or the initial state leads approximately to the same belief, i.e., having reached the target.

Thank you!