

Inria
INVENTEURS DU MONDE NUMÉRIQUE

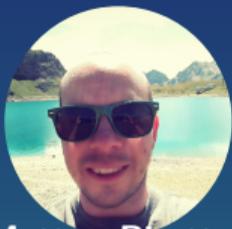


Artificial Intelligence Research

Near Optimal Exploration-Exploitation in Non-Communicating Markov Decision Processes



Ronan Fruit[†]



Matteo Pirotta^{*}



Alessandro Lazaric^{*}

[†]Sequel – INRIA Lille

^{*}FAIR – Facebook Paris

Inria
INVENTEURS DU MONDE NUMÉRIQUE

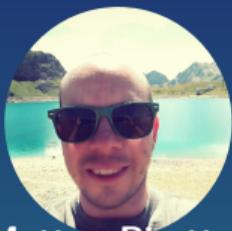


facebook
Artificial Intelligence Research

Exploration–exploitation in RL with Misspecified State Space



Ronan Fruit[†]



Matteo Pirota*



Alessandro Lazaric*

[†]Sequel – INRIA Lille

*FAIR – Facebook Paris

Misspecified states: Examples

1 Breakout [Mnih et al., 2015]



Misspecified states: Examples

1 Breakout [Mnih et al., 2015]

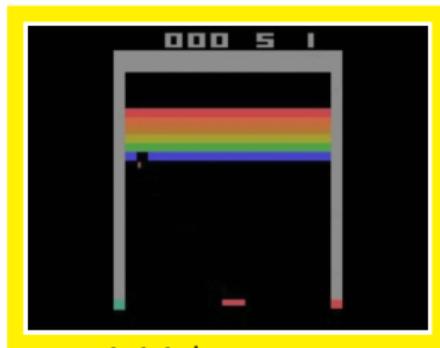
Intuitive state space: set of plausible configurations of wall, ball and paddle



Misspecified states: Examples

1 Breakout [Mnih et al., 2015]

Intuitive state space: set of plausible configurations of wall, ball and paddle



initial state s_1



Misspecified states: Examples

1 Breakout [Mnih et al., 2015]

Intuitive state space: set of plausible configurations of wall, ball and paddle



initial state s_1



Plausible state after some time...



Misspecified states: Examples

1 Breakout [Mnih et al., 2015]

Intuitive state space: set of plausible configurations of wall, ball and paddle



initial state s_1



Plausible state after some time...



Non reachable from s_1

Misspecified states: Examples

1 Breakout [Mnih et al., 2015]

Intuitive state space: set of plausible configurations of wall, ball and paddle



initial state s_1



Plausible state after some time...

Cannot be observed!

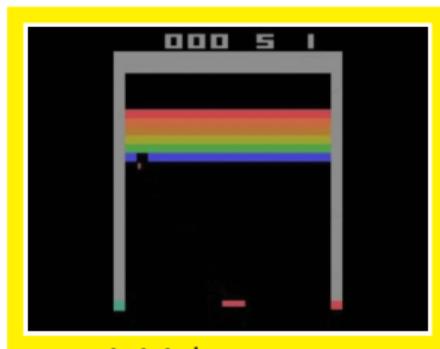


Non reachable from s_1

Misspecified states: Examples

1 Breakout [Mnih et al., 2015]

Intuitive state space: set of plausible configurations of wall, ball and paddle



initial state s_1



Plausible state after some time...

Cannot be observed!



Non reachable from s_1

Misspecified state space = \exists states non-observable from initial state
 + difficult to exclude **explicitly** from the state space

Why is exploration more challenging with a misspecified state space?

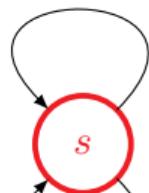
Why is exploration more challenging with a misspecified state space?

- All existing methods known to efficiently balance exploration and exploitation in RL with theoretical guarantees rely on the optimism in the face of uncertainty principle
- All such methods fail to learn when the state space is misspecified

Why is exploration more challenging with a misspecified state space?

- All existing methods known to **efficiently** balance exploration and exploitation in RL with **theoretical guarantees** rely on the **optimism in the face of uncertainty** principle
- All such methods **fail to learn** when the state space is **misspecified**

$$a_0, r_0 = 0$$



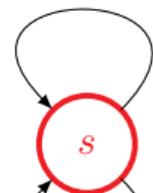
$$a_1, r_1 = \frac{1}{2}$$

Example 1 of Ortner [2008]

Why is exploration more challenging with a misspecified state space?

- All existing methods known to efficiently balance exploration and exploitation in RL with theoretical guarantees rely on the optimism in the face of uncertainty principle
- All such methods **fail to learn** when the state space is **misspecified**

$$a_0, r_0 = 0$$



$$a_1, r_1 = \frac{1}{2}$$

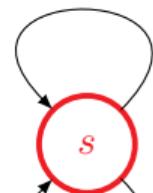
Optimism (UCB, etc.) = Optimal Strategy

Example 1 of Ortner [2008]

Why is exploration more challenging with a misspecified state space?

- All existing methods known to **efficiently** balance exploration and exploitation in RL with **theoretical guarantees** rely on the **optimism in the face of uncertainty** principle
- All such methods **fail to learn** when the state space is **misspecified**

$$a_0, r_0 = 0$$



$$a_1, r_1 = \frac{1}{2}$$

Not reachable from s

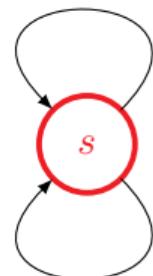


Example 1 of Ortner [2008]

Why is exploration more challenging with a misspecified state space?

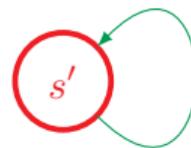
- All existing methods known to **efficiently** balance exploration and exploitation in RL with **theoretical guarantees** rely on the **optimism in the face of uncertainty** principle
- All such methods **fail to learn** when the state space is **misspecified**

$$a_0, r_0 = 0$$



$$a_1, r_1 = \frac{1}{2}$$

Not reachable from s



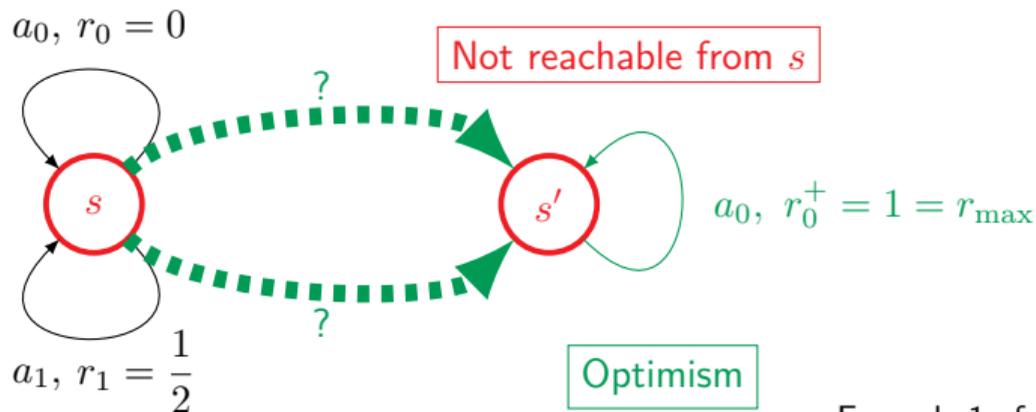
$$a_0, r_0^+ = 1 = r_{\max}$$

Optimism

Example 1 of Ortner [2008]

Why is exploration more challenging with a misspecified state space?

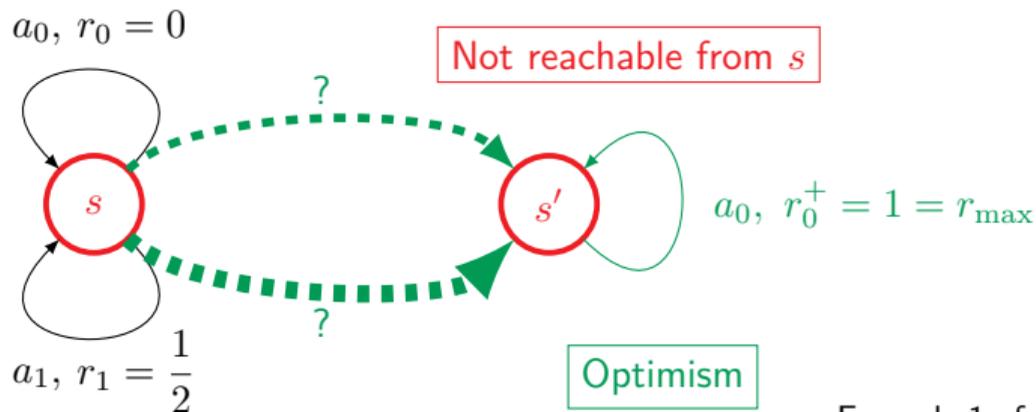
- All existing methods known to **efficiently** balance exploration and exploitation in RL with **theoretical guarantees** rely on the **optimism in the face of uncertainty** principle
- All such methods **fail to learn** when the state space is **misspecified**



Example 1 of Ortner [2008]

Why is exploration more challenging with a misspecified state space?

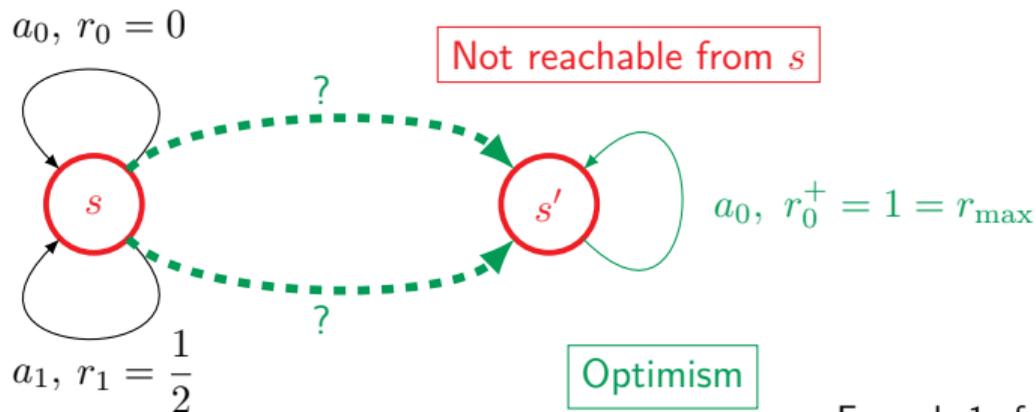
- All existing methods known to **efficiently** balance exploration and exploitation in RL with **theoretical guarantees** rely on the **optimism in the face of uncertainty** principle
- All such methods **fail to learn** when the state space is **misspecified**



Example 1 of Ortner [2008]

Why is exploration more challenging with a misspecified state space?

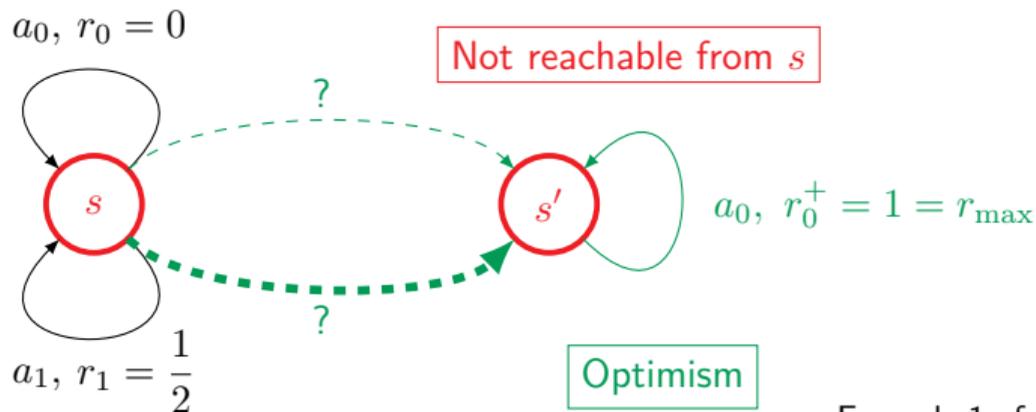
- All existing methods known to **efficiently** balance exploration and exploitation in RL with **theoretical guarantees** rely on the **optimism in the face of uncertainty** principle
- All such methods **fail to learn** when the state space is **misspecified**



Example 1 of Ortner [2008]

Why is exploration more challenging with a misspecified state space?

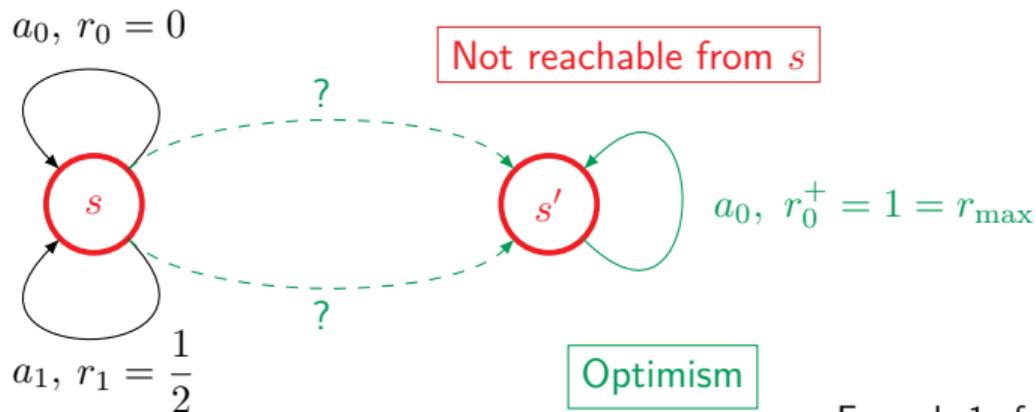
- All existing methods known to **efficiently** balance exploration and exploitation in RL with **theoretical guarantees** rely on the **optimism in the face of uncertainty** principle
- All such methods **fail to learn** when the state space is **misspecified**



Example 1 of Ortner [2008]

Why is exploration more challenging with a misspecified state space?

- All existing methods known to **efficiently** balance exploration and exploitation in RL with **theoretical guarantees** rely on the **optimism in the face of uncertainty** principle
- All such methods **fail to learn** when the state space is **misspecified**

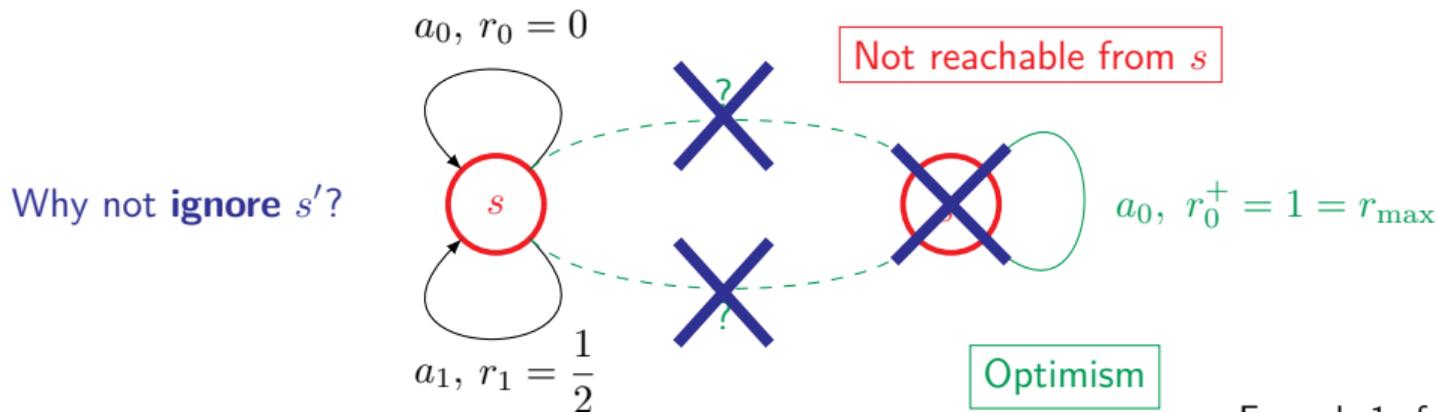


Example 1 of Ortner [2008]

Problem: The action played keeps changing: it is a_0 half of the time and a_1 the other half \implies **linear regret!**

Why is exploration more challenging with a misspecified state space?

- All existing methods known to **efficiently** balance exploration and exploitation in RL with **theoretical guarantees** rely on the **optimism in the face of uncertainty** principle
- All such methods **fail to learn** when the state space is **misspecified**

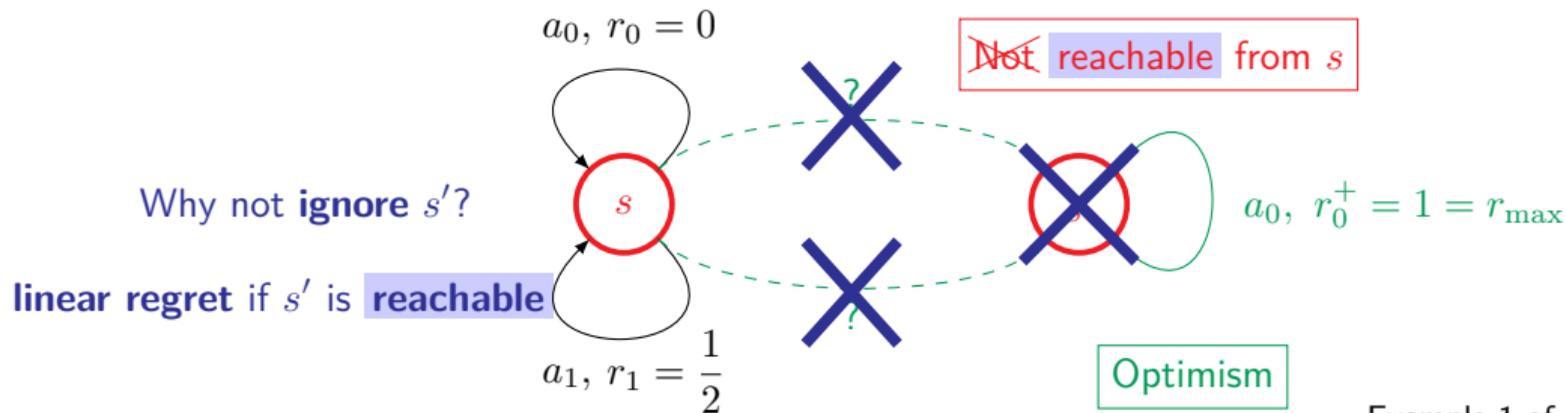


Example 1 of Ortner [2008]

Problem: The action played keeps changing: it is a_0 half of the time and a_1 the other half \implies **linear regret!**

Why is exploration more challenging with a misspecified state space?

- All existing methods known to efficiently balance exploration and exploitation in RL with theoretical guarantees rely on the optimism in the face of uncertainty principle
- All such methods **fail to learn** when the state space is **misspecified**



Example 1 of Ortner [2008]

Problem: The action played keeps changing: it is a_0 half of the time and a_1 the other half \Rightarrow **linear regret!**

Our work

👉 **Regret** of existing methods: $\tilde{O}\left(D S \sqrt{AT}\right)$




Diameter Total number of states

👉 **Misspecified** state space $\iff D = +\infty$ (infinite diameter)

👉 TUCRL: first algorithm able to **adapt** to the **reachable part** of the MDP

Regret:

$\tilde{O}\left(D^{\text{C}} S^{\text{C}} \sqrt{AT}\right)$




Reachable diameter Number of reachable states

Come to see our poster # 161 !