

# BEYOND OPTIMISM: EXPLORATION WITH PARTIALLY OBSERVABLE REWARDS

**SIMONE PARISI**  
ALIREZA KAZEMIPOUR  
MICHAEL BOWLING



38TH CONFERENCE ON NEURAL INFORMATION PROCESSING SYSTEMS (NEURIPS 2024)

# OPTIMISM WORKS WHEN REWARDS ARE OBSERVABLE



- Small and easy-to-find reward to the left.
- Large and hard-to-find reward to the far right.
- Optimistic algorithms (e.g., Q-Learning with optimistic initialization) would solve this problem easily.

# BUT **FAILS** IF REWARDS ARE PARTIALLY OBSERVABLE



- Coin rewards are observable only if the agent **pushes the button first**.
- But the agent will also receives **negative rewards** in the meantime.

## BUT **FAILS** IF REWARDS ARE PARTIALLY OBSERVABLE



- Coin rewards are observable only if the agent **pushes the button first**.
- But the agent will also receives **negative rewards** in the meantime.
- Coin rewards **do exist** even when not observed, so the optimal policy is to ignore the button and collect the large coin ...
- ... but the agent must push the button first (**suboptimal action**) to learn the optimal policy!

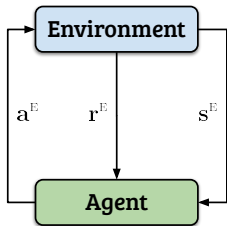
## BUT **FAILS** IF REWARDS ARE PARTIALLY OBSERVABLE



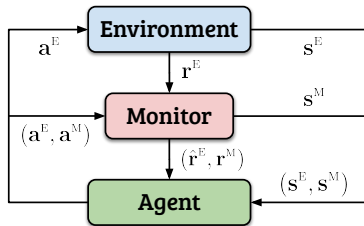
- Coin rewards are observable only if the agent **pushes the button first**.
- But the agent will also receives **negative rewards** in the meantime.
- Coin rewards **do exist** even when not observed, so the optimal policy is to ignore the button and collect the large coin ...
- ... but the agent must push the button first (**suboptimal action**) to learn the optimal policy!

***Optimism  $\neq$  Suboptimal Actions***

# MONITORED MDPs



MDP



Monitored MDP

- The **monitor is a separate MDP** “on the top” of the environment.
- The goal is to maximize  $\sum_{t=1}^{\infty} \gamma^t (r_t^E + r_t^M) \dots$
- ... but the agent **observes proxy rewards**  $\hat{r}_t^E \sim \mathcal{M}(r_t^E, s_t^M, a_t^M)$  instead of true rewards  $r_t^E \dots$
- ... and **sometimes**  $\hat{r}_t^E = \text{NaN!}$

[“Monitored Markov Decision Processes”, AAMAS 2024]

# DIRECTED REWARD-FREE EXPLORATION-EXPLOITATION

```
1  $(s^G, a^G) = \arg \min_{s,a} N_t(s, a)$  // goal: the least-visited state-action pair
2  $\beta_t = \log t / N_t(s^G, a^G)$  // how much did we visit it?
3 if  $\beta_t > \bar{\beta}$  then return  $\rho(a | s_t, s^G, a^G)$  // explore: follow goal-conditioned policy
4 else return  $\arg \max_a \hat{Q}(s_t, a)$  // exploit: follow the greedy policy
```

# DIRECTED REWARD-FREE EXPLORATION-EXPLOITATION

```
1  $(s^G, a^G) = \arg \min_{s,a} N_t(s, a)$  // goal: the least-visited state-action pair
2  $\beta_t = \log t / N_t(s^G, a^G)$  // how much did we visit it?
3 if  $\beta_t > \bar{\beta}$  then return  $\rho(a | s_t, s^G, a^G)$  // explore: follow goal-conditioned policy
4 else return  $\arg \max_a Q(s_t, a)$  // exploit: follow the greedy policy
```



# DIRECTED REWARD-FREE EXPLORATION-EXPLOITATION

```
1  $(s^G, a^G) = \arg \min_{s,a} N_t(s, a)$  // goal: the least-visited state-action pair
2  $\beta_t = \log t / N_t(s^G, a^G)$  // how much did we visit it?
3 if  $\beta_t > \bar{\beta}$  then return  $\rho(a | s_t, s^G, a^G)$  // explore: follow goal-conditioned policy
4 else return  $\arg \max_a Q(s_t, a)$  // exploit: follow the greedy policy
```

$\rho(a | s_t, s^G, a^G)$  should **reach any goal pair “as fast as possible”**. How?

# DIRECTED REWARD-FREE EXPLORATION-EXPLOITATION

- 1  $(s^G, a^G) = \arg \min_{s,a} N_t(s, a)$  // goal: the least-visited state-action pair
- 2  $\beta_t = \log t / N_t(s^G, a^G)$  // how much did we visit it?
- 3 if  $\beta_t > \bar{\beta}$  then return  $\rho(a | s_t, s^G, a^G)$  // explore: follow goal-conditioned policy
- 4 else return  $\arg \max_a Q(s_t, a)$  // exploit: follow the greedy policy

$\rho(a | s_t, s^G, a^G)$  should **reach any goal pair “as fast as possible”**. How?  
It **maximizes the goal occurrences**, i.e.,

$$\rho^*(s | a, s_i, a_j) = \arg \max_{\rho} S_{s_i a_j}^{\rho}(s, a)$$

$$S_{s_i a_j}^{\rho}(s_t, a_t) = \mathbb{E} \left[ \sum_{k=t}^{\infty} \gamma^{k-t} \mathbb{1}_{\{s_k=s_i, a_k=a_j\}} \right] \quad \Leftarrow \text{S-function}$$

$$\mathbb{1}_{\{s_k=s_i, a_k=a_j\}} = \begin{cases} 1 & \text{if } s_k = s_i \text{ and } a_k = a_j \\ 0 & \text{otherwise} \end{cases} \quad \Leftarrow \text{Successor features}$$

# DIRECTED REWARD-FREE EXPLORATION-EXPLOITATION

- 1  $(s^G, a^G) = \arg \min_{s,a} N_t(s, a)$  // goal: the least-visited state-action pair
- 2  $\beta_t = \log t / N_t(s^G, a^G)$  // how much did we visit it?
- 3 **if**  $\beta_t > \bar{\beta}$  **then return**  $\rho(a | s_t, s^G, a^G)$  // explore: follow goal-conditioned policy
- 4 **else return**  $\arg \max_a Q(s_t, a)$  // exploit: follow the greedy policy

$\rho(a | s_t, s^G, a^G)$  should **reach any goal pair “as fast as possible”**. How?  
It **maximizes the goal occurrences**, i.e.,

$$\rho^*(s | a, s_i, a_j) = \arg \max_{\rho} S_{s_i a_j}^{\rho}(s, a)$$

$$S_{s_i a_j}^{\rho}(s_t, a_t) = \mathbb{E} \left[ \sum_{k=t}^{\infty} \gamma^{k-t} \mathbb{1}_{\{s_k=s_i, a_k=a_j\}} \right] \quad \Leftarrow \text{S-function}$$

$$\mathbb{1}_{\{s_k=s_i, a_k=a_j\}} = \begin{cases} 1 & \text{if } s_k = s_i \text{ and } a_k = a_j \\ 0 & \text{otherwise} \end{cases} \quad \Leftarrow \text{Successor features}$$

**Exploration does not depend on the reward observability!**

# RESULTS (EXPECTED RETURN)

