

Unsupervised Reinforcement Learning

Theoretical Guarantees in the Hard and Easy Cases

Jingfeng Wu, Vladimir Braverman, Lin F. Yang

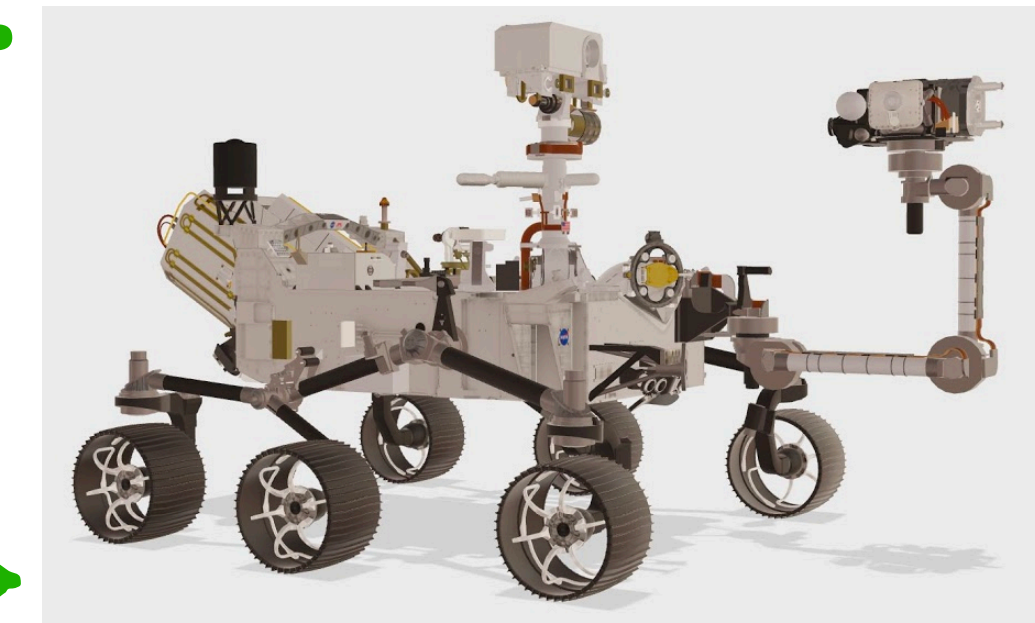
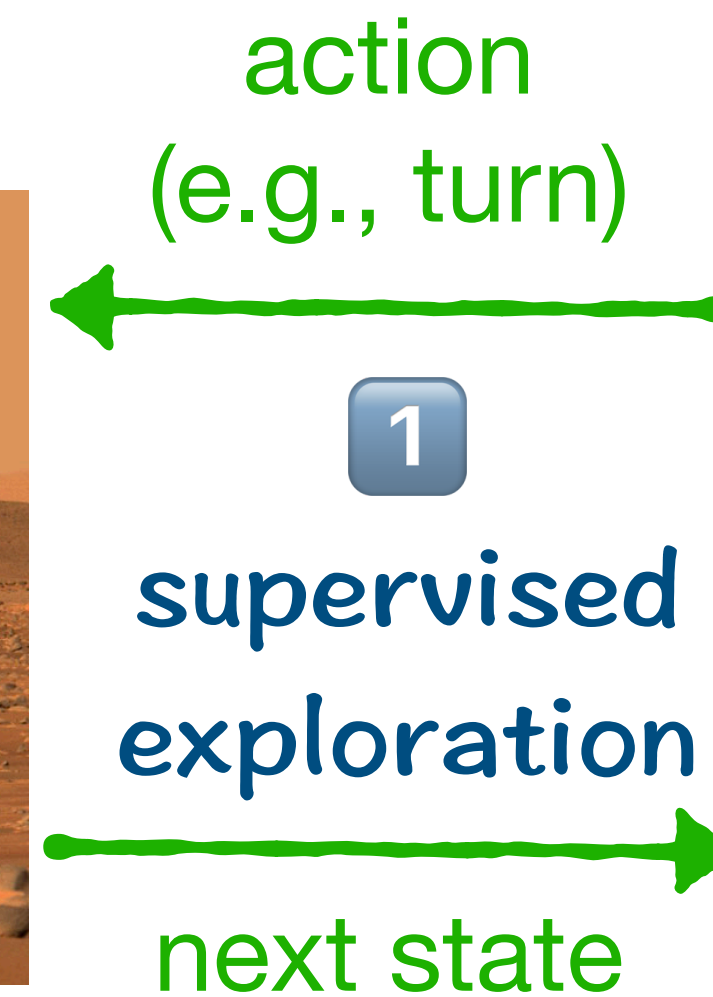


Supervised RL

- states: \mathcal{S}
- actions: \mathcal{A}
- horizon length: H
- transition kernel: $\mathbf{P} : \mathcal{S} \times \mathcal{A} \rightarrow \mu(\mathcal{S})$
- reward function:
 $r : [H] \times \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$
- policy:
 $\pi : \mathcal{S} \rightarrow \mu(\mathcal{A})$



Mars



Perseverance Rover

(weather, location, etc..) and reward (water?)

2 compute a “good” policy π



$$\mathbb{P}\{V_1^* - V_1^\pi > \epsilon\} < \delta$$

$$Q_h^\pi(x, a) := \mathbb{E}\left[r_h(x_h, a_h) + \dots + r_H(x_H, a_H)\right]$$

$$V_h^\pi(x) := Q_h^\pi(x, \pi_h(x))$$

$$Q^* := \max_{\pi} Q^\pi$$

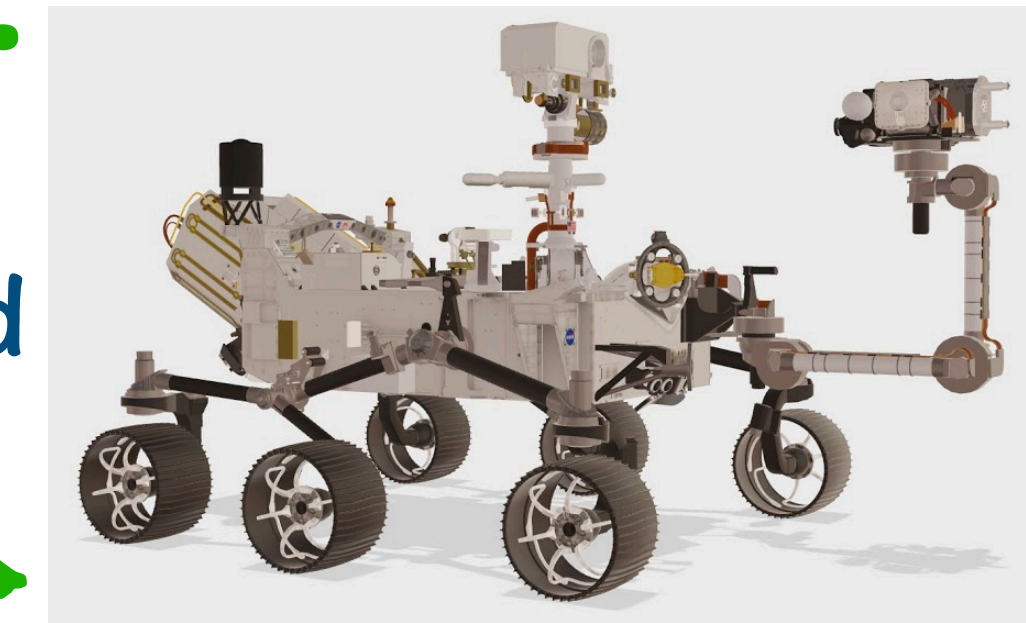
$$V^* := \max_{\pi} V^\pi$$

Unsupervised RL

\approx MDP +
a set of reward functions



Mars



Perseverance Rover

$$\mathcal{R} \subset \{r : [H] \times \mathcal{S} \times \mathcal{A} \rightarrow [0,1]\}$$

2 user chooses a
task from \mathcal{R}

“arbitrary”

for an “independent”

$r \in \mathcal{R}$

$\mathbb{P} \left\{ V_1^*(r) - V_1^\pi(r) > \epsilon \right\} < \delta$

Reward-Free Exploration

Task-Agnostic Exploration



3 compute a “good” policy π



Jin, C., Krishnamurthy, A., Simchowitz, M., & Yu, T. (2020, November). Reward-free exploration for reinforcement learning. In *International Conference on Machine Learning* (pp. 4870-4879). PMLR.

Zhang, X., Ma, Y., & Singla, A. (2020). Task-agnostic exploration in reinforcement learning. *Advances in Neural Information Processing Systems*, 33, 11734-11743.

Reduction to Supervised RL

[Algorithm] For each $r \in \mathcal{R}$, learning a policy π with a supervised RL algorithm.

$$\mathcal{R} \subset \{r : [H] \times \mathcal{S} \times \mathcal{A} \rightarrow [0,1]\}$$

[Sample Complexity]

$$K = \mathcal{O}(|\mathcal{R}| \cdot H^2 SA \cdot \log / \epsilon^2)$$

[Memory] $\propto |\mathcal{R}|$, costly

$$\approx \frac{1}{\epsilon^{SA}}$$

[Supervised RL] K trajectories are sufficient/necessary to solve supervised RL

$$K = \Theta(H^2 SA \cdot \log / \epsilon^2)$$

Azar, M. G., Osband, I., & Munos, R. (2017, July). Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning* (pp. 263-272). PMLR.

Dann, C., & Brunskill, E. (2015). Sample complexity of episodic fixed-horizon reinforcement learning. *Advances in Neural Information Processing Systems*, 28.

$\hat{\mathbf{P}}$ + Dynamic Programming

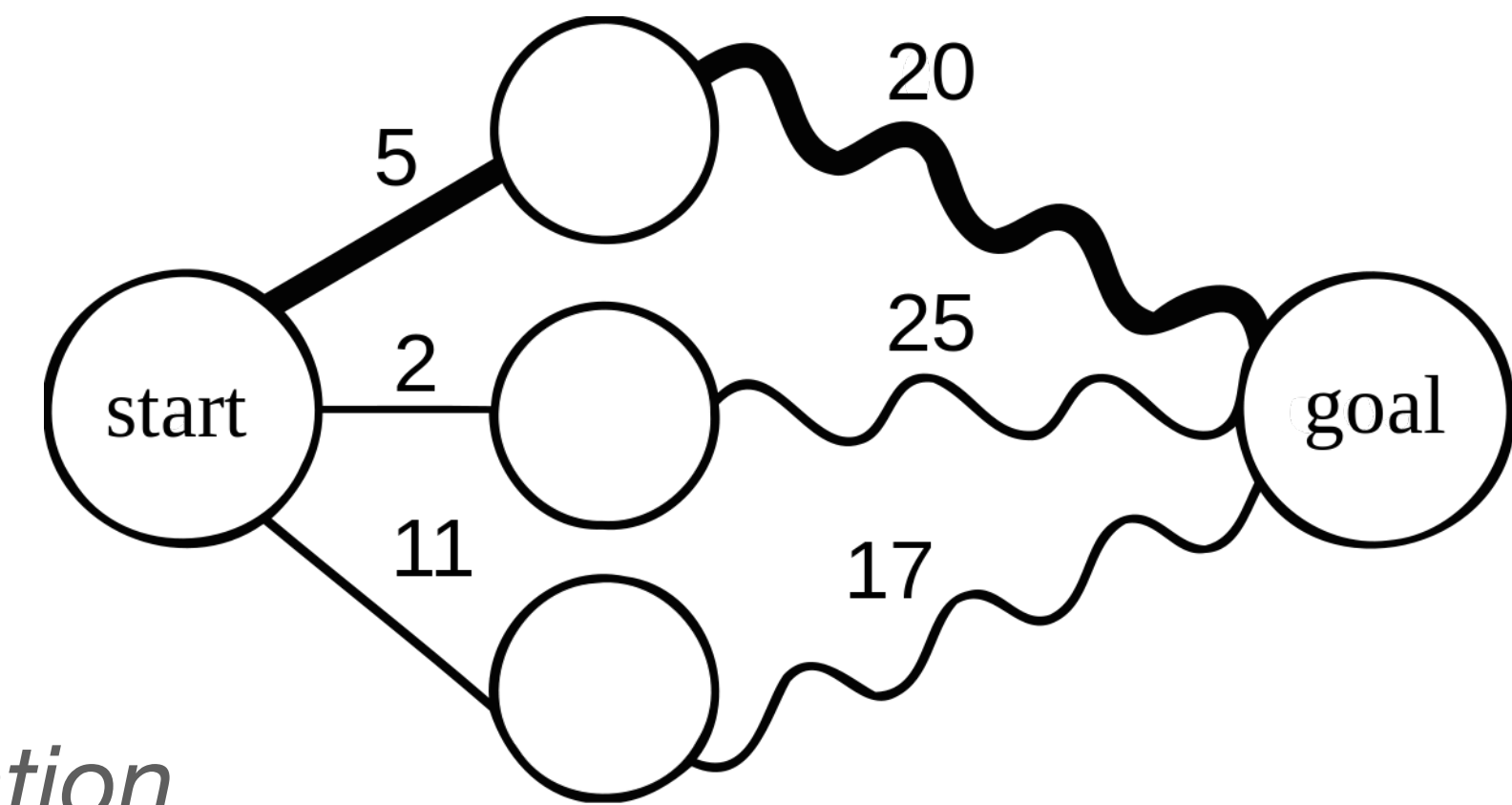
[Algorithm] Estimating $\hat{\mathbf{P}} \approx \mathbf{P} \in [0,1]^{S^2A}$,

then DP for each r

[Sample Complexity] $K \propto S^2A / \epsilon^2$

[How $\hat{\mathbf{P}}$?] w/ generative model 😊,

otherwise 🤔 [Jin et. al, 2020]



Bellman Equation

$$Q_h^*(x, a) := r_h(x, a) + \max_{a \in \mathcal{A}} \mathbb{E}_{y \sim \mathbf{P}(\cdot|x, a)} V_{h+1}^*(y)$$

$$V_h^*(x) := \max_{a \in \mathcal{A}} Q_h^*(x, a)$$

Necessary to have an accurate model?

S^2 factor: yes for RFE, no for TAE

[Total Variation] With N i.i.d samples, $|\hat{\mathbf{P}} - \mathbf{P}|_{\ell_1} < \sqrt{S^2A \cdot \log / N}$, w.h.p.

Minimax Cases $\mathcal{R} := \{r : [H] \times \mathcal{S} \times \mathcal{A} \rightarrow [0,1]\}$

[Question] how much exploration (K) is sufficient/necessary to compute (ϵ, δ) -correct policy?

$$\mathbb{P} \left\{ \begin{array}{l} \text{for an "independent"} \\ r \in \mathcal{R} \\ V_1^*(r) - V_1^\pi(r) > \epsilon \end{array} \right\} < \delta$$

[Upper Bound] There is an ALO that needs at most K trajectories:

$$K = \mathcal{O}(H^3 SA \cdot \log / \epsilon^2)$$

[Lower Bound] Every $(\epsilon, 0.1)$ -correct ALO needs at least K trajectories:

$$\mathbb{E}[K] \geq \Omega(H^2 SA / \epsilon^2)$$

Unsupervised RL is nearly as hard as supervised RL

Wu, J., Braverman, V., & Yang, L. (2021). Accommodating picky customers: Regret bound and exploration complexity for multi-objective reinforcement learning. *Advances in Neural Information Processing Systems*, 34.

Dann, C., & Brunskill, E. (2015). Sample complexity of episodic fixed-horizon reinforcement learning. *Advances in Neural Information Processing Systems*, 28.

UCBVI

Algorithm 1 UCBVI

```

Initialize data  $\mathcal{H} = \emptyset$ 
for episode  $k = 1, 2, \dots, K$  do
     $Q_{k,h} = \text{UCB-Q-values}(\mathcal{H})$ 
    for step  $h = 1, \dots, H$  do
        Take action  $a_{k,h} = \arg \max_a Q_{k,h}(x_{k,h}, a)$ 
        Update  $\mathcal{H} = \mathcal{H} \cup (x_{k,h}, a_{k,h}, x_{k,h+1})$ 
    end for
end for

```

[Analysis]

$$\begin{aligned}
 V_1^* - V_1^\pi &\lesssim \bar{V}_1^\pi \\
 &\leq \bar{V}_1^K \\
 &\leq \frac{1}{K} \cdot \sum_{k=1}^K \bar{V}_1^k \\
 &\lesssim \frac{1}{K} \cdot \sqrt{H^3 S A K}
 \end{aligned}$$

$$R(x, a) = 0$$

$$\text{bonus}^k(x, a) \approx \sqrt{\frac{H^2 \log}{N^k(x, a)}} + \text{lower orders}$$

“reward-independent bonus”

Algorithm 2 UCB-Q-values

Require: Bonus algorithm bonus, Data \mathcal{H}

```

Compute, for all  $(x, a, y) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$ ,
 $N_k(x, a, y) = \sum_{(x', a', y') \in \mathcal{H}} \mathbb{I}(x' = x, a' = a, y' = y)$ 
 $N_k(x, a) = \sum_{y \in \mathcal{S}} N_k(x, a, y)$ 
 $N'_{k,h}(x, a) = \sum_{(x_{i,h}, a_{i,h}, x_{i,h+1}) \in \mathcal{H}} \mathbb{I}(x_{i,h} = x, a_{i,h} = a)$ 
Let  $\mathcal{K} = \{(x, a) \in \mathcal{S} \times \mathcal{A}, N_k(x, a) > 0\}$ 
Estimate  $\hat{P}_k(y|x, a) = \frac{N_k(x, a, y)}{N_k(x, a)}$  for all  $(x, a) \in \mathcal{K}$ 
Initialize  $V_{k,H+1}(x) = 0$  for all  $(x, a) \in \mathcal{S} \times \mathcal{A}$ 
for  $h = H, H-1, \dots, 1$  do
    for  $(x, a) \in \mathcal{S} \times \mathcal{A}$  do
        if  $(x, a) \in \mathcal{K}$  then
             $b_{k,h}(x, a) = \text{bonus}(\hat{P}_k, V_{k,h+1}, N_k, N'_{k,h})$ 
             $Q_{k,h}(x, a) = \min(Q_{k-1,h}(x, a), H,$ 
                 $R(x, a) + (\hat{P}_k V_{k,h+1})(x, a) + b_{k,h}(x, a))$ 
        else
             $Q_{k,h}(x, a) = H$ 
        end if
         $V_{k,h}(x) = \max_{a \in \mathcal{A}} Q_{k,h}(x, a)$ 
    end for
end for
return Q-values  $Q_{k,h}$ 

```

Azar, M. G., Osband, I., & Munos, R. (2017, July). Minimax regret bounds for reinforcement learning. In *International Conference on Machine Learning* (pp. 263-272). PMLR.

Wu, J., Braverman, V., & Yang, L. (2021). Accommodating picky customers: Regret bound and exploration complexity for multi-objective reinforcement learning. *Advances in Neural Information Processing Systems*, 34.

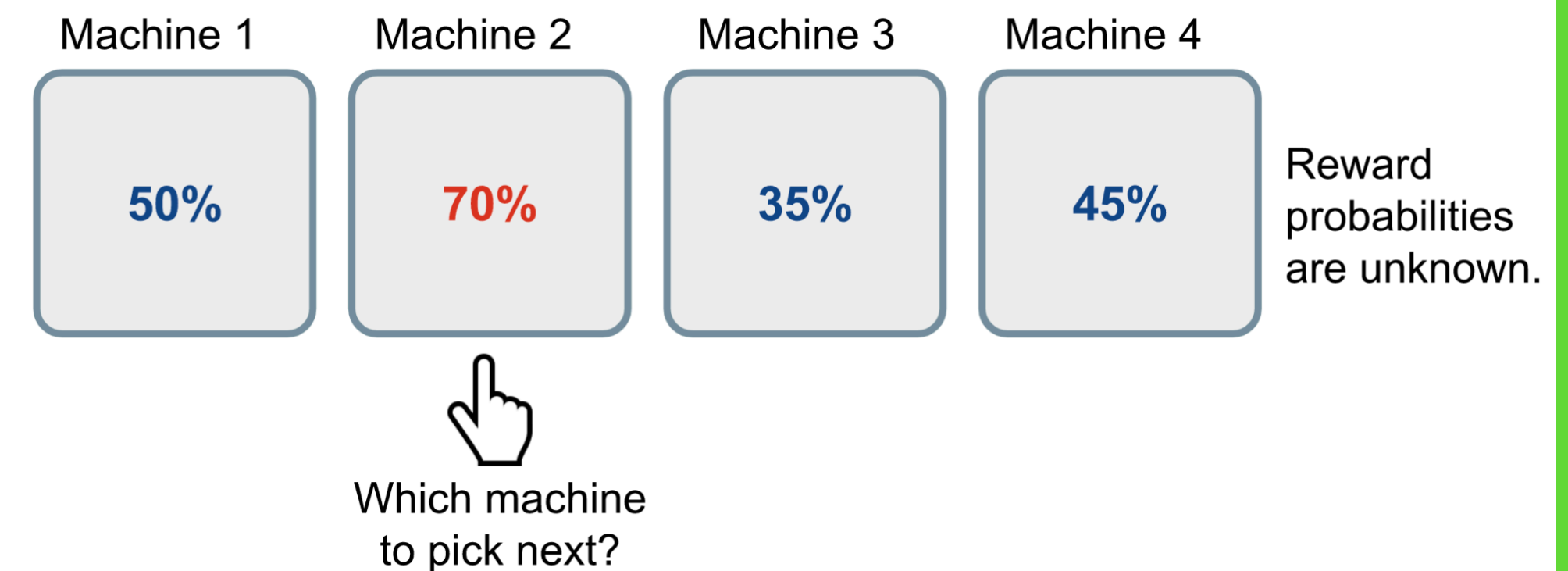
Gap Cases $\mathcal{R} := \{r : \text{gap}(r) \geq \rho\}$

$$\text{gap}(r) := \min_{x,a,h} \{\text{nonzero } V_h^*(x; r) - Q_h^*(x, a; r)\}$$

[Sample Complexity] $\approx \begin{cases} \tilde{\Theta}(1), & H = 1 \\ ?, & H \geq 2 \end{cases}$

*For unsupervised bandits ($H=1$),
gap enables an acceleration $\tilde{\Theta}(1/\epsilon^2) \rightarrow \tilde{\Theta}(1)$*

[Unsupervised Bandit ($H=1$)]



A Bandit Picture from [Lil'Log](#)

[Algorithm] Uniform exploration

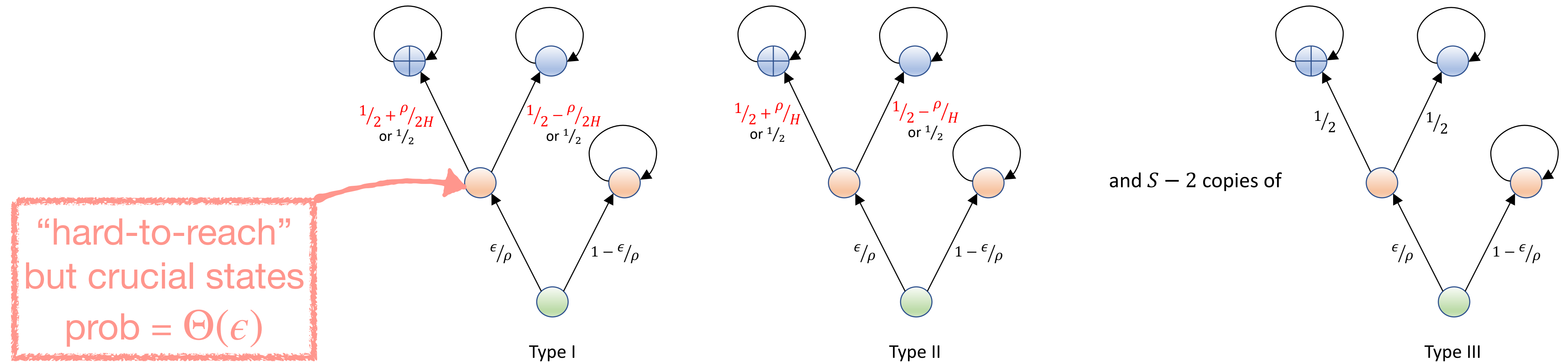
[Analysis]

$$\begin{aligned} \mathbb{P}\{a' \neq a\} &= \mathbb{P}\left\{\hat{R}_{a'} > \hat{R}_a\right\} \\ &= \mathbb{P}\left\{\left(\hat{R}_{a'} - r_{a'}\right) - \left(\hat{R}_a - r_a\right) > r_a - r_{a'}\right\} \\ &\leq \mathbb{P}\left\{\left(\hat{R}_{a'} - r_{a'}\right) - \left(\hat{R}_a - r_a\right) > \rho\right\} \\ &\lesssim A \exp(-\rho^2 K) \approx A \exp(-\rho^2 T/A). \end{aligned}$$

A Lower Bound

[Lower Bound] Any (ϵ, δ) -correct ALO in gap cases needs at least K episodes,

$$\mathbb{E}[K] \geq \begin{cases} \Omega\left(\frac{H^2 SA}{\rho \epsilon} \cdot \log \frac{1}{\delta}\right) = \Omega\left(\frac{1}{\epsilon}\right), & H \geq 2; \\ \Omega\left(\frac{SA}{\rho^2} \log \frac{1}{\delta}\right) = \Omega(1), & H = 1. \end{cases}$$



An Algorithm and An Upper Bound

[Exploration] “Modified UCBVI”

- “reward” $\rightarrow 0$
- bonus is *clipped* (set to zero if it is small) (ρ is an input)

$$\text{bonus}^k(x, a) \approx \text{clip}_{\frac{\rho}{H}} \left(\sqrt{\frac{H^2 \log}{N^k(x, a)}} \right) + \text{lower orders}$$

[Planning] The usual UCBVI

$$\text{bonus}^k(x, a) \approx \sqrt{\frac{H^2 \log}{N^k(x, a)}}$$

[Upper Bound] There is an (ϵ, δ) -correct ALO, that needs K episodes

$$K \leq \tilde{\mathcal{O}} \left(\frac{H^3 S A}{\rho \epsilon} \cdot \log \frac{1}{\delta} + \frac{H^4 S^2 A}{\epsilon} \cdot \log \frac{1}{\delta} \right) = \tilde{\mathcal{O}} \left(\frac{1}{\epsilon} \right)$$

where $\tilde{\mathcal{O}}$ hides $\log^2(HSAK)$ and constants.

For unsupervised RL, gap enables an acceleration $\tilde{\Theta}(1/\epsilon^2) \rightarrow \tilde{\Theta}(1/\epsilon)$

Take-Home

1. Unsupervised RL \approx supervised RL

$$\text{unsupervised} \propto \tilde{\mathcal{O}}(H^3 SA/\epsilon^2)$$

vs.

$$\text{Supervised} \propto \tilde{\mathcal{O}}(H^2 SA/\epsilon^2)$$

2. gap-cases are easier, but is still “hard” when $H \geq 2$

$$\text{gap-rate} \propto \tilde{\mathcal{O}}(1/\epsilon)$$

vs.

$$\text{minimax-rate} \propto \tilde{\mathcal{O}}(1/\epsilon^2)$$

Open Problems

1. Improving H dependence?
2. An algorithm agnostic to ρ ?
3. Removing lower order S^2 ?

