

### Authors

Sebastian Reboul  
Hélène Halconruy

### Domains

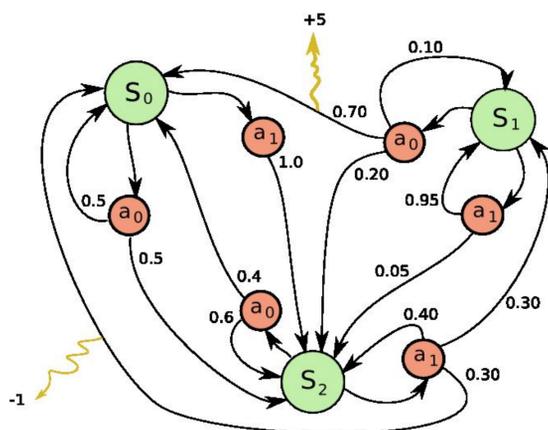
Reinforcement Learning  
Probability Theory

### Key Words

Regret Bounds  
Concentration Inequalities

### Partners

Wiremind Cargo



### Setting

**Markov Decision Process (MDP)** learning Algorithm.

**Objective** : Produce a sequence of policies that minimize cumulative regret.

### Regret

$$V_h^\pi(s) = \mathbb{E}_\pi \left[ \sum_{t=h}^H r_t(s_t, a_t) \mid s_h = s \right]$$

$$\text{Regret}(T) = \sum_{t=1}^T \left( V_1^*(s_1^t) - V_1^{\pi^t}(s_1^t) \right)$$

### Online RL

1. Interactive
2. Can adapt strategy on the fly
3. Does not have outside help from another player

### Offline RL

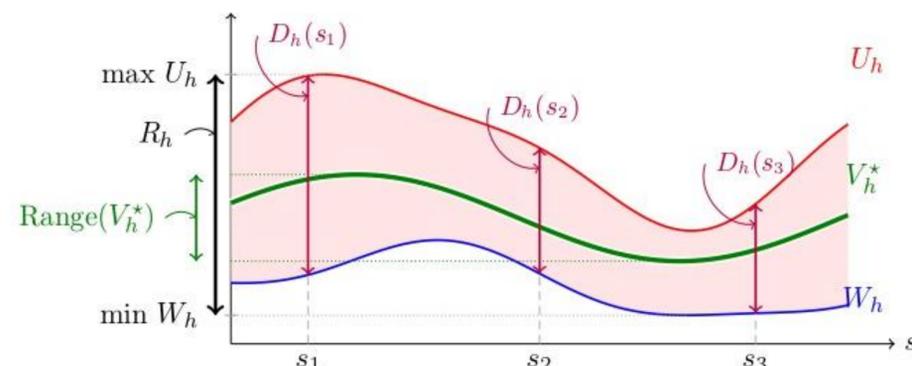
1. Non Interactive
2. Access to games from a single player
3. Works great if that player is good

### Offline-to-Online

1. First get data from **any** player (K games)
2. Extract information (privacy preserving)
3. Use it to accelerate online learning.

### Framework

Run an offline algorithm to produce **Value Envelopes** (high probability bounds on the value function)



**Strategy**: Pick actions maximising :

$$b_h^{t,\text{on}}(s, a) = c_1 \sigma_{h+1}^t(s, a) \sqrt{\frac{L}{N_h^t(s, a)}} + c_2 \frac{R_{h+1} L}{N_h^t(s, a)}$$

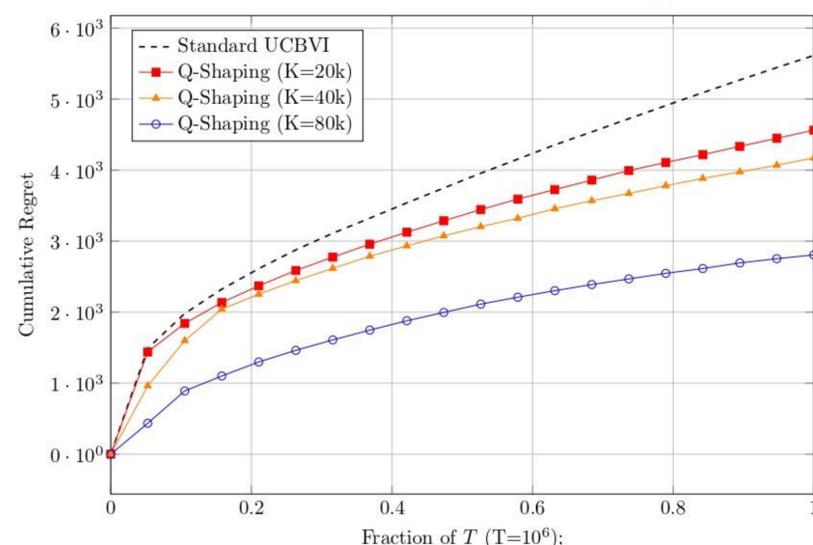
where :

$$\sigma_{h+1}^t(s, a) := \sqrt{\text{Var}_{s' \sim \hat{P}_t(\cdot|s, a)} \left( \frac{1}{2} (U_{h+1}(s') + W_{h+1}(s')) \right) + \frac{1}{2} \sqrt{\mathbb{E}_{s' \sim \hat{P}_t(\cdot|s, a)} [D_{h+1}(s')^2]}}$$

### Results

$$\begin{aligned} \text{Regret}(T) \leq & \tilde{O} \left( R^{\max} \left( \sqrt{TH} |\text{PairEff}| + |\text{PairEff}| + \sqrt{T} \right) \right. \\ & + \sqrt{\frac{H^5}{K d_{\min}^b}} \left( \sqrt{TH} |\text{PairEff}| + |S|H |\text{PairEff}| + |S|H \sqrt{T} \right) \\ & \left. + \frac{H^3}{(K d_{\min}^b)} \left( \sqrt{TH} |\text{PairEff}| + |S|H |\text{PairEff}| + |S|H \sqrt{T} \right) \right) \end{aligned}$$

Q-Shaping Performance vs. Offline Data Size (K)



Link to Preprint

