



POLITECNICO
MILANO 1863

Exploration in Policy Search by Multiple Importance Sampling

Lorenzo Lupo
lorenzo.lupo@mail.polimi.it

April 16th, 2019

How can robots learn to backflip?



Algorithmic Trading



Self-driving Cars



Plan

1. Basics of Reinforcement Learning
2. Exploration in Policy Search
3. Problem Formalization
4. OPTIMIST

The Reinforcement Learning Framework

The Reinforcement Learning Framework

Environment



\mathcal{P}, \mathcal{R}

The Reinforcement Learning Framework

Agent



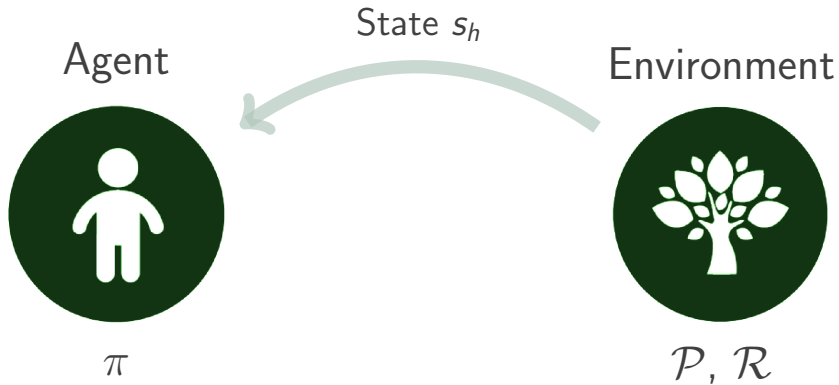
π

Environment

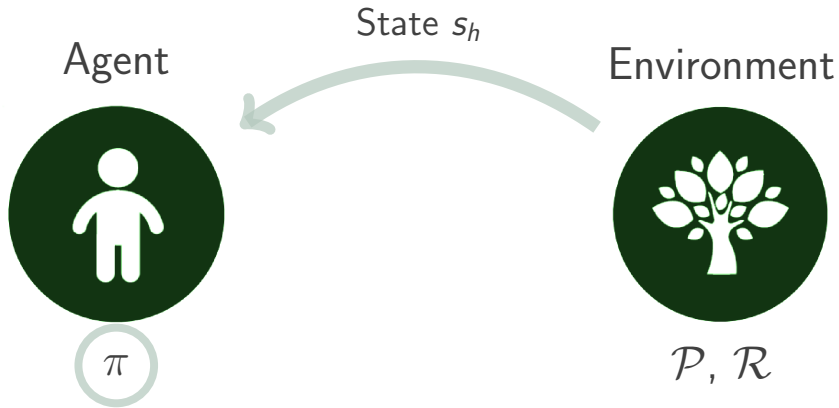


\mathcal{P}, \mathcal{R}

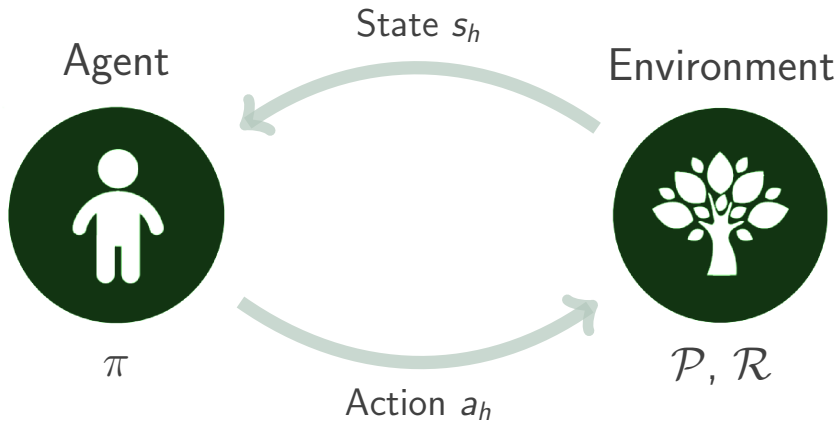
The Reinforcement Learning Framework



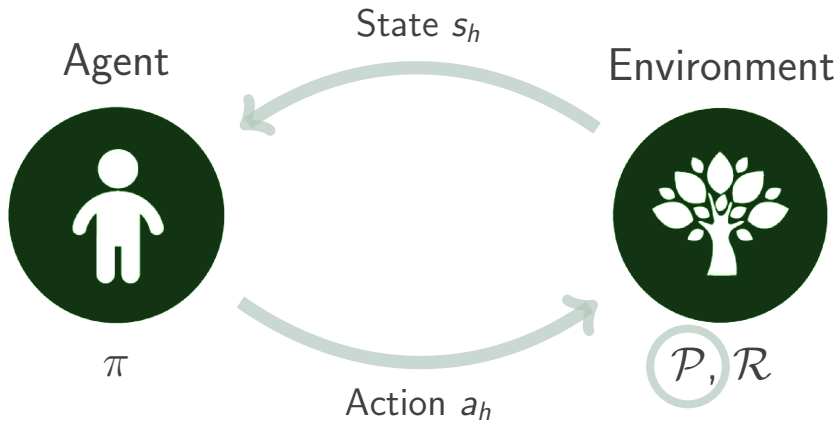
The Reinforcement Learning Framework



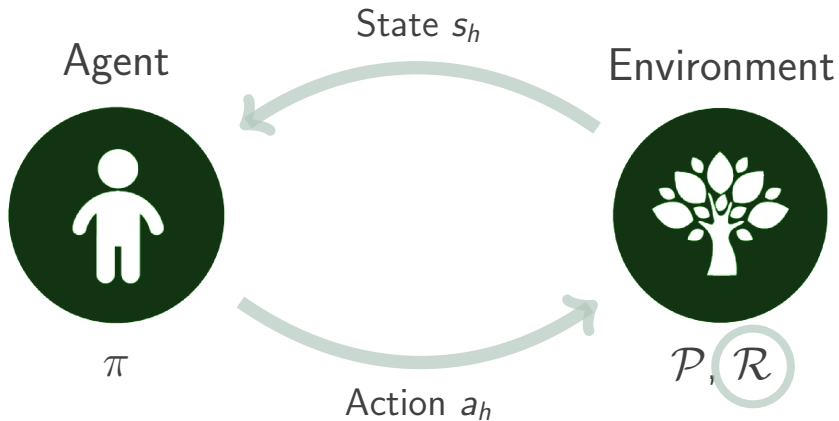
The Reinforcement Learning Framework



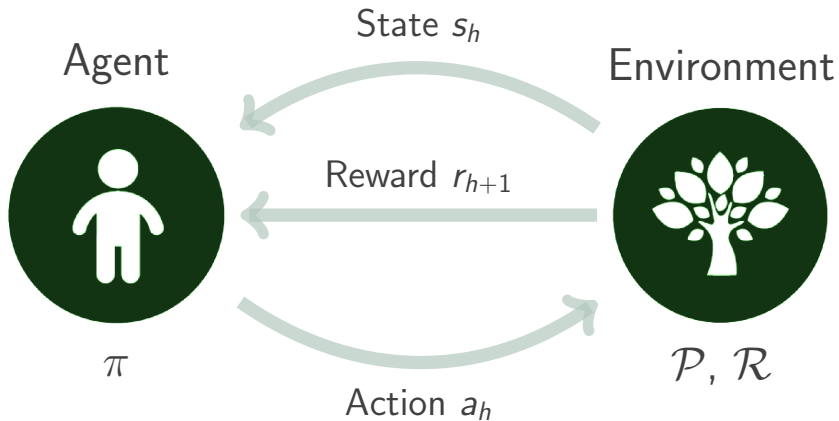
The Reinforcement Learning Framework



The Reinforcement Learning Framework



The Reinforcement Learning Framework



Policy Search Formulation

Cumulative return of a trajectory τ :

$$\mathcal{R}(\tau) = \sum_{h=0}^{H-1} \gamma^h r_{h+1}, \text{ with } \tau = [s_0, a_0, s_1, a_1, \dots, s_{H-1}, a_{H-1}]$$

Policy Search Formulation

Cumulative return of a trajectory τ :

$$\mathcal{R}(\tau) = \sum_{h=0}^{H-1} \gamma^h r_{h+1}, \text{ with } \tau = [s_0, a_0, s_1, a_1, \dots, s_{H-1}, a_{H-1}]$$

Parametric policy:

$$\pi_{\theta} : \mathcal{S} \rightarrow \Delta(\mathcal{A}), \text{ i.e., } \pi_{\theta}(a|s) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{1}{2} \left(\frac{a - \theta^T \phi(s)}{\sigma} \right)^2 \right)$$

Policy Search Formulation

Cumulative return of a trajectory τ :

$$\mathcal{R}(\tau) = \sum_{h=0}^{H-1} \gamma^h r_{h+1}, \text{ with } \tau = [s_0, a_0, s_1, a_1, \dots, s_{H-1}, a_{H-1}]$$

Parametric policy:

$$\pi_{\theta} : \mathcal{S} \rightarrow \Delta(\mathcal{A}), \text{ i.e., } \pi_{\theta}(a|s) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{1}{2} \left(\frac{a - \theta^T \phi(s)}{\sigma} \right)^2 \right)$$

Performance:

$\mu(\theta) = \mathbb{E}_{\tau \sim p_{\theta}} [\mathcal{R}(\tau)]$, where p_{θ} is the **distribution over trajectories** $\tau \in \mathcal{T}$ induced by π_{θ}

Policy Search Formulation

Cumulative return of a trajectory τ :

$$\mathcal{R}(\tau) = \sum_{h=0}^{H-1} \gamma^h r_{h+1}, \text{ with } \tau = [s_0, a_0, s_1, a_1, \dots, s_{H-1}, a_{H-1}]$$

Parametric policy:

$$\pi_{\theta} : \mathcal{S} \rightarrow \Delta(\mathcal{A}), \text{ i.e., } \pi_{\theta}(a|s) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{1}{2} \left(\frac{a - \theta^T \phi(s)}{\sigma} \right)^2 \right)$$

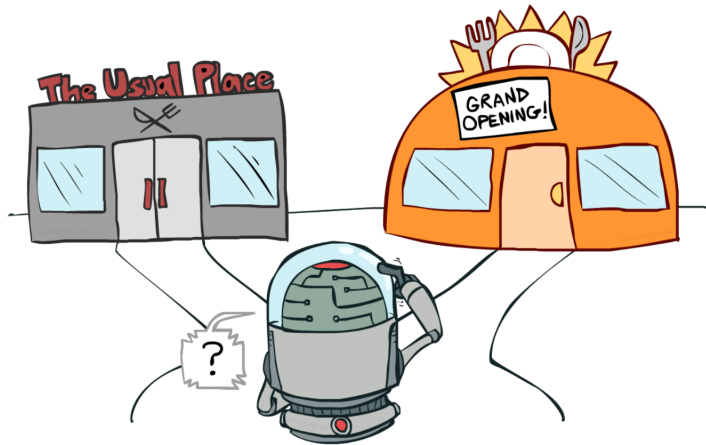
Performance:

$\mu(\theta) = \mathbb{E}_{\tau \sim p_{\theta}} [\mathcal{R}(\tau)]$, where p_{θ} is the **distribution over trajectories** $\tau \in \mathcal{T}$ induced by π_{θ}

Objective:

$$\theta^* = \arg \max_{\theta \in \Theta} \mu(\theta).$$

Exploration VS Exploitation



Plan

1. Basics of Reinforcement Learning
2. Exploration in Policy Search
3. Problem Formalization
4. OPTIMIST

Exploration in Policy Search

Undirected exploration

Generate actions based on randomness, without any knowledge of the learning process.

Exploration in Policy Search

Undirected exploration

Generate actions based on randomness, without any knowledge of the learning process.

- Ex1: by adopting stochastic policies.

Exploration in Policy Search

Undirected exploration

Generate actions based on randomness, without any knowledge of the learning process.

- Ex1: by adopting stochastic policies.
- Ex2: by augmenting rewards with the entropy of the policy:

$$\mathcal{R}(\tau) = \sum_{h=0}^{H-1} \gamma^h r_{h+1} + \mathcal{H}(\pi_{\theta}(\cdot|s_h)).$$

Exploration in Policy Search

Undirected exploration

Generate actions based on randomness, without any knowledge of the learning process.

- Ex1: by adopting stochastic policies.
- Ex2: by augmenting rewards with the entropy of the policy:

$$\mathcal{R}(\tau) = \sum_{h=0}^{H-1} \gamma^h r_{h+1} + \mathcal{H}(\pi_{\theta}(\cdot | s_h)).$$

Directed exploration

Leverage on the knowledge acquired during learning.

Exploration in Policy Search

Undirected exploration

Generate actions based on randomness, without any knowledge of the learning process.

- Ex1: by adopting stochastic policies.
- Ex2: by augmenting rewards with the entropy of the policy:

$$\mathcal{R}(\tau) = \sum_{h=0}^{H-1} \gamma^h r_{h+1} + \mathcal{H}(\pi_{\theta}(\cdot | s_h)).$$

Directed exploration

Leverage on the knowledge acquired during learning.

- Count-based techniques.

Plan

1. Basics of Reinforcement Learning
2. Exploration in Policy Search
- 3. Problem Formalization**
4. OPTIMIST

Problem Formalization

Decision Set or Arms Set

The parameter space $\Theta \subseteq \mathbb{R}^d$.

Problem Formalization

Decision Set or Arms Set

The parameter space $\Theta \subseteq \mathbb{R}^d$.

Procedure

At every decision step $t \in [0, 1, 2, \dots, T]$:

Problem Formalization

Decision Set or Arms Set

The parameter space $\Theta \subseteq \mathbb{R}^d$.

Procedure

At every decision step $t \in [0, 1, 2, \dots, T]$:

1. **Select** an arm $\theta_t \in \Theta$;

Problem Formalization

Decision Set or Arms Set

The parameter space $\Theta \subseteq \mathbb{R}^d$.

Procedure

At every decision step $t \in [0, 1, 2, \dots, T]$:

1. **Select** an arm $\theta_t \in \Theta$;
2. **Sample** a trajectory $\tau_t \in \mathcal{T}$ by following π_{θ_t} ;

Problem Formalization

Decision Set or Arms Set

The parameter space $\Theta \subseteq \mathbb{R}^d$.

Procedure

At every decision step $t \in [0, 1, 2, \dots, T]$:

1. **Select** an arm $\theta_t \in \Theta$;
2. **Sample** a trajectory $\tau_t \in \mathcal{T}$ by following π_{θ_t} ;
3. **Observe** the cumulative return $\mathcal{R}(\tau_t)$.

Problem Formalization

Decision Set or Arms Set

The parameter space $\Theta \subseteq \mathbb{R}^d$.

Procedure

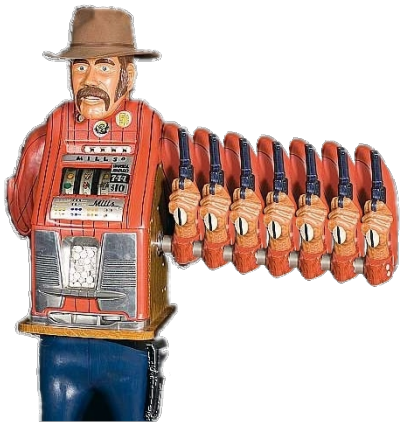
At every decision step $t \in [0, 1, 2, \dots, T]$:

1. **Select** an arm $\theta_t \in \Theta$;
2. **Sample** a trajectory $\tau_t \in \mathcal{T}$ by following π_{θ_t} ;
3. **Observe** the cumulative return $\mathcal{R}(\tau_t)$.

Goal

Minimize $\text{Regret}(T) = \sum_{t=0}^T \mu(\theta^*) - \mu(\theta_t)$, where $\theta^* = \arg \max_{\theta \in \Theta} \mu(\theta)$

Problem Formulation



Multi Armed Bandits

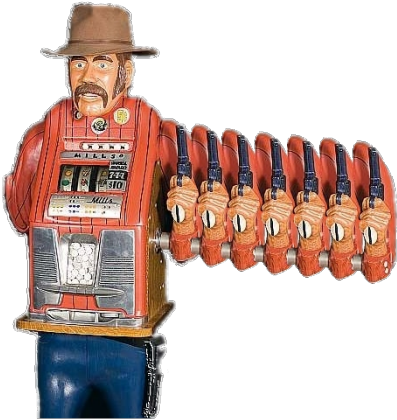
Problem Formulation



Multi Armed Bandits

- Simpler framework;

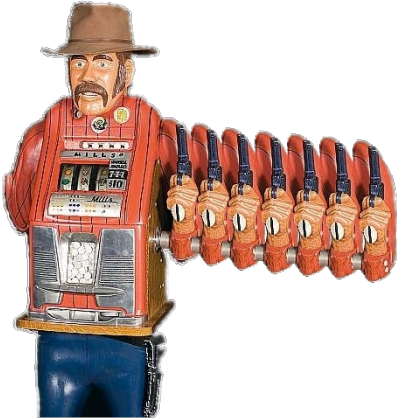
Problem Formulation



Multi Armed Bandits

- Simpler framework;
- Share the exploration-exploitation tradeoff;

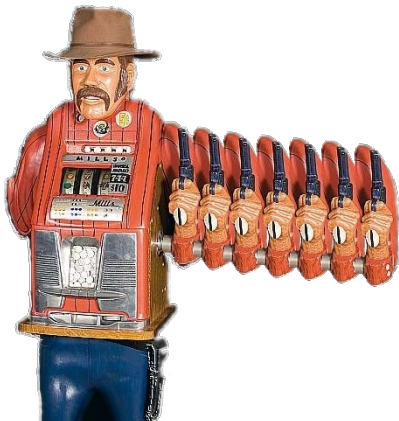
Problem Formulation



Multi Armed Bandits

- Simpler framework;
- Share the exploration-exploitation tradeoff;
- Ample literature available;

Problem Formulation



Multi Armed Bandits

- Simpler framework;
- Share the exploration-exploitation tradeoff;
- Ample literature available;

Desideratum

sub-linear $\text{Regret}(T) \Leftrightarrow \lim_{T \rightarrow \infty} \text{Regret}(T)/T = 0$

E.g. $\text{Regret}(T) = \mathcal{O}(\log T)$

Plan

1. Basics of Reinforcement Learning
2. Exploration in Policy Search
3. Problem Formalization
4. OPTIMIST

Optimistic Policy Optimization via MIS with Truncation

Algorithm 1 OPTIMIST

1: **Input:** initial parameters θ_0 , confidence schedule $(\delta_t)_{t=1}^T$, order $\epsilon \in (0, 1]$

Optimistic Policy Optimization via MIS with Truncation

Algorithm 2 OPTIMIST

- 1: **Input:** initial parameters θ_0 , confidence schedule $(\delta_t)_{t=1}^T$, order $\epsilon \in (0, 1]$
- 2: Draw trajectory $\tau_0 \sim p_{\theta_0}$ and observe return $\mathcal{R}(\tau_0)$

Optimistic Policy Optimization via MIS with Truncation

Algorithm 3 OPTIMIST

- 1: **Input:** initial parameters θ_0 , confidence schedule $(\delta_t)_{t=1}^T$, order $\epsilon \in (0, 1]$
- 2: Draw trajectory $\tau_0 \sim p_{\theta_0}$ and observe return $\mathcal{R}(\tau_0)$
- 3: **for** $t = 1, \dots, T$ **do**

Optimistic Policy Optimization via MIS with Truncation

Algorithm 4 OPTIMIST

- 1: **Input:** initial parameters θ_0 , confidence schedule $(\delta_t)_{t=1}^T$, order $\epsilon \in (0, 1]$
- 2: Draw trajectory $\tau_0 \sim p_{\theta_0}$ and observe return $\mathcal{R}(\tau_0)$
- 3: **for** $t = 1, \dots, T$ **do**
- 4: Select arm $\theta_t = \arg \max_{\theta \in \Theta} B_t^\epsilon(\theta, \delta_t)$

Optimistic Policy Optimization via MIS with Truncation

Algorithm 5 OPTIMIST

- 1: **Input:** initial parameters θ_0 , confidence schedule $(\delta_t)_{t=1}^T$, order $\epsilon \in (0, 1]$
 - 2: Draw trajectory $\tau_0 \sim p_{\theta_0}$ and observe return $\mathcal{R}(\tau_0)$
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: Select arm $\theta_t = \arg \max_{\theta \in \Theta} B_t^\epsilon(\theta, \delta_t)$
 - 5: Draw trajectory $\tau_t \sim p_{\theta_t}$ and observe return $\mathcal{R}(\tau_t)$
 - 6: **end for**
-

Optimistic Policy Optimization via MIS with Truncation

Algorithm 6 OPTIMIST

- 1: **Input:** initial parameters θ_0 , confidence schedule $(\delta_t)_{t=1}^T$, order $\epsilon \in (0, 1]$
 - 2: Draw trajectory $\tau_0 \sim p_{\theta_0}$ and observe return $\mathcal{R}(\tau_0)$
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: Select arm $\theta_t = \arg \max_{\theta \in \Theta} B_t^\epsilon(\theta, \delta_t)$
 - 5: Draw trajectory $\tau_t \sim p_{\theta_t}$ and observe return $\mathcal{R}(\tau_t)$
 - 6: **end for**
-

$$\mu(\theta) \leq B_t^\epsilon(\theta, \delta_t) = \check{\mu}_t^{MIS}(\theta) + \|f\|_\infty \left(\sqrt{2} + \frac{4}{3} \right) \left(\frac{d_{1+\epsilon}(p_{\theta_t} \|\Phi_t) \log \frac{1}{\delta_t}}{t} \right)^{\frac{\epsilon}{1+\epsilon}}$$

Optimistic Policy Optimization via MIS with Truncation

Algorithm 7 OPTIMIST

- 1: **Input:** initial parameters θ_0 , confidence schedule $(\delta_t)_{t=1}^T$, order $\epsilon \in (0, 1]$
 - 2: Draw trajectory $\tau_0 \sim p_{\theta_0}$ and observe return $\mathcal{R}(\tau_0)$
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: Select arm $\theta_t = \arg \max_{\theta \in \Theta} B_t^\epsilon(\theta, \delta_t)$
 - 5: Draw trajectory $\tau_t \sim p_{\theta_t}$ and observe return $\mathcal{R}(\tau_t)$
 - 6: **end for**
-

Upper Confidence Bound

$$\boxed{\mu(\theta) \leq B_t^\epsilon(\theta, \delta_t)} = \check{\mu}_t^{MIS}(\theta) + \|f\|_\infty \left(\sqrt{2} + \frac{4}{3} \right) \left(\frac{d_{1+\epsilon}(p_{\theta_t} \|\Phi_t) \log \frac{1}{\delta_t}}{t} \right)^{\frac{\epsilon}{1+\epsilon}}$$

Optimistic Policy Optimization via MIS with Truncation

Algorithm 8 OPTIMIST

- 1: **Input:** initial parameters θ_0 , confidence schedule $(\delta_t)_{t=1}^T$, order $\epsilon \in (0, 1]$
 - 2: Draw trajectory $\tau_0 \sim p_{\theta_0}$ and observe return $\mathcal{R}(\tau_0)$
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: Select arm $\theta_t = \arg \max_{\theta \in \Theta} B_t^\epsilon(\theta, \delta_t)$
 - 5: Draw trajectory $\tau_t \sim p_{\theta_t}$ and observe return $\mathcal{R}(\tau_t)$
 - 6: **end for**
-

Truncated Multiple Importance Sampling Estimator

$$\mu(\theta) \leq B_t^\epsilon(\theta, \delta_t) = \boxed{\check{\mu}_t^{MIS}(\theta)} + \|f\|_\infty \left(\sqrt{2} + \frac{4}{3} \right) \left(\frac{d_{1+\epsilon}(p_{\theta_t} \|\Phi_t) \log \frac{1}{\delta_t}}{t} \right)^{\frac{\epsilon}{1+\epsilon}}$$

Optimistic Policy Optimization via MIS with Truncation

Algorithm 9 OPTIMIST

- 1: **Input:** initial parameters θ_0 , confidence schedule $(\delta_t)_{t=1}^T$, order $\epsilon \in (0, 1]$
 - 2: Draw trajectory $\tau_0 \sim p_{\theta_0}$ and observe return $\mathcal{R}(\tau_0)$
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: Select arm $\theta_t = \arg \max_{\theta \in \Theta} B_t^\epsilon(\theta, \delta_t)$
 - 5: Draw trajectory $\tau_t \sim p_{\theta_t}$ and observe return $\mathcal{R}(\tau_t)$
 - 6: **end for**
-

Exploration Bonus

$$\mu(\theta) \leq B_t^\epsilon(\theta, \delta_t) = \check{\mu}_t^{MIS}(\theta) + \left\| f \right\|_\infty \left(\sqrt{2} + \frac{4}{3} \right) \left(\frac{d_{1+\epsilon}(p_{\theta_t} \| \Phi_t) \log \frac{1}{\delta_t}}{t} \right)^{\frac{\epsilon}{1+\epsilon}}$$

Optimistic Policy Optimization via MIS with Truncation

Algorithm 10 OPTIMIST






- 1: **Input:** initial parameters θ_0 , confidence schedule $(\delta_t)_{t=1}^T$, order $\epsilon \in (0, 1]$
 - 2: Draw trajectory $\tau_0 \sim p_{\theta_0}$ and observe return $\mathcal{R}(\tau_0)$
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: Select arm $\theta_t = \arg \max_{\theta \in \Theta} B_t^\epsilon(\theta, \delta_t)$
 - 5: Draw trajectory $\tau_t \sim p_{\theta_t}$ and observe return $\mathcal{R}(\tau_t)$
 - 6: **end for**
-

Exploration Bonus

$$\mu(\theta) \leq B_t^\epsilon(\theta, \delta_t) = \check{\mu}_t^{MIS}(\theta) + \|f\|_\infty \left(\sqrt{2} + \frac{4}{3} \right) \left(\frac{d_{1+\epsilon}(p_{\theta_t} \|\Phi_t) \log \frac{1}{\delta_t}}{t} \right)^{\frac{\epsilon}{1+\epsilon}}$$

Thank you for your attention!

References

-  Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. (2011).
Improved algorithms for linear stochastic bandits.
In Advances in Neural Information Processing Systems, pages 2312–2320.
-  Agrawal, R. (1995a).
The continuum-armed bandit problem.
SIAM Journal on Control and Optimization, 33(6):1926–1951.
-  Agrawal, R. (1995b).
Sample mean based index policies by $o(\log n)$ regret for the multi-armed bandit problem.
Advances in Applied Probability, 27(4):1054–1078.
-  Agrawal, S. and Goyal, N. (2013).
Further optimal regret bounds for thompson sampling.
In Artificial intelligence and statistics, pages 99–107.
-  Amari, S.-I. (1998).