



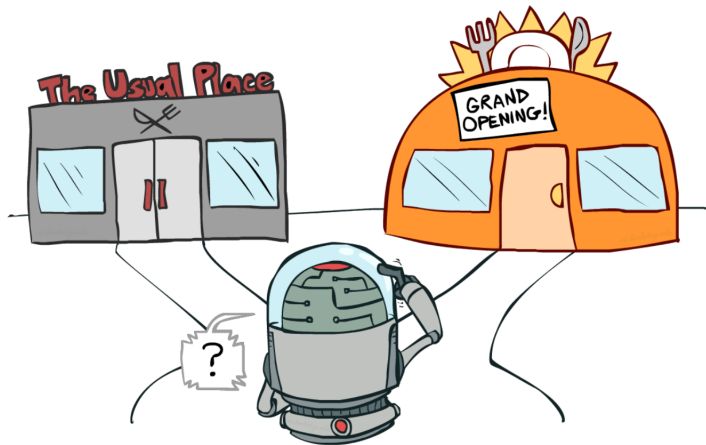
POLITECNICO
MILANO 1863

Exploration in Policy Search by Multiple Importance Sampling

Lorenzo Lupo
lorenzo.lupo@mail.polimi.it

April 16th, 2019

Exploration VS Exploitation



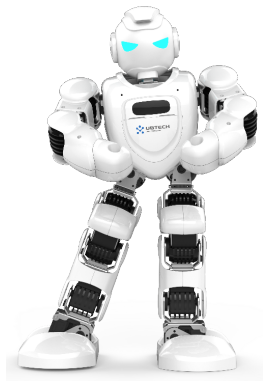
Reinforcement Learning Applications



Reinforcement Learning Applications



Reinforcement Learning Applications



Plan

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

The Reinforcement Learning Framework

The Reinforcement Learning Framework

Environment



\mathcal{P}, \mathcal{R}

The Reinforcement Learning Framework

Agent



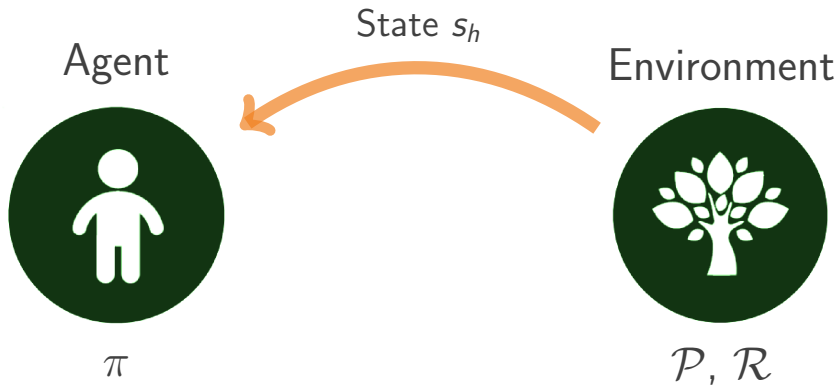
π

Environment

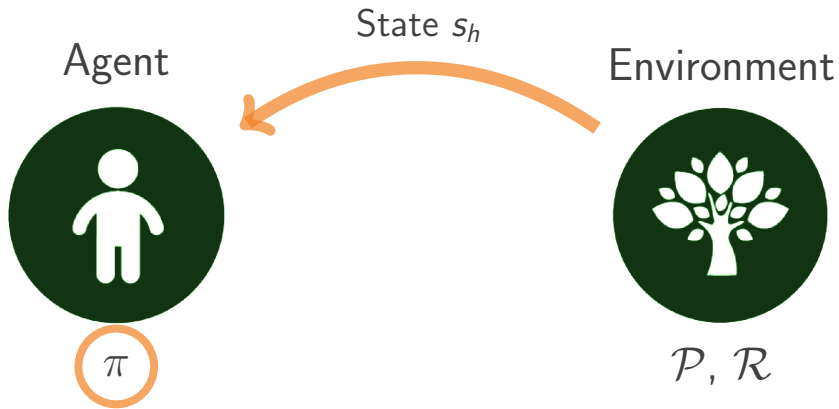


\mathcal{P}, \mathcal{R}

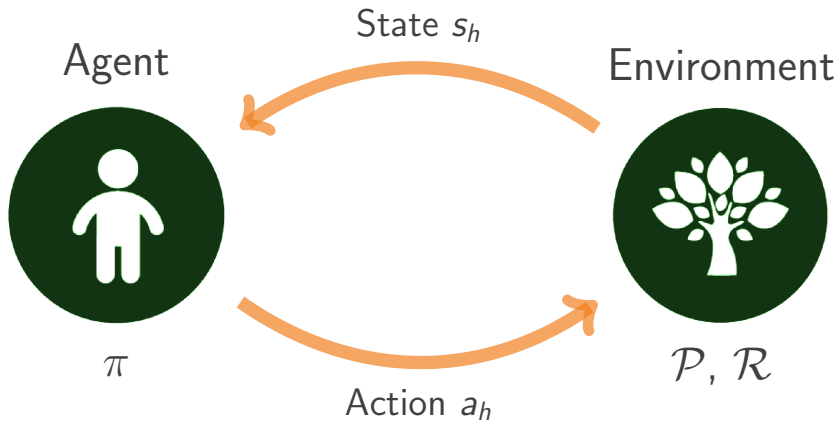
The Reinforcement Learning Framework



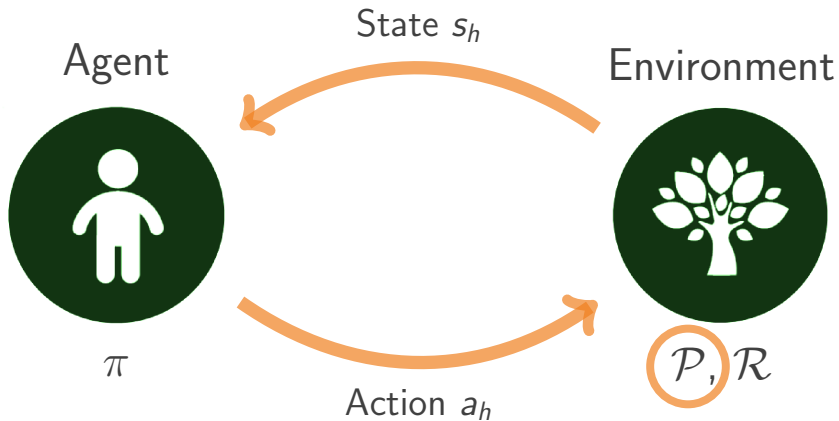
The Reinforcement Learning Framework



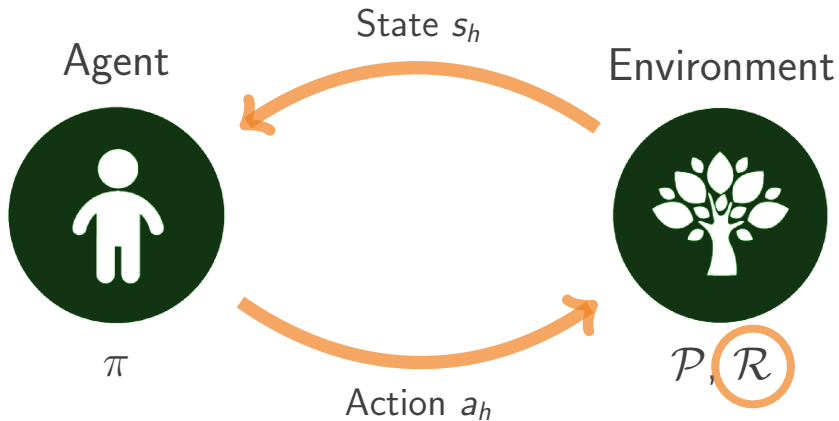
The Reinforcement Learning Framework



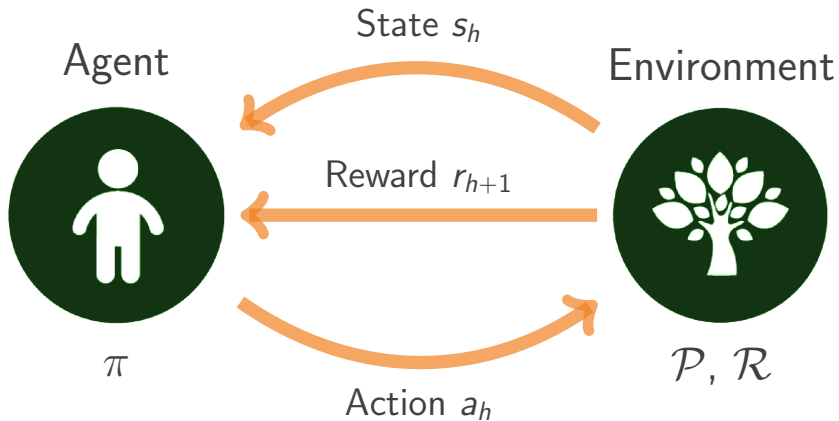
The Reinforcement Learning Framework



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:

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

Plan

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

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
5. Experiments
6. Conclusions

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

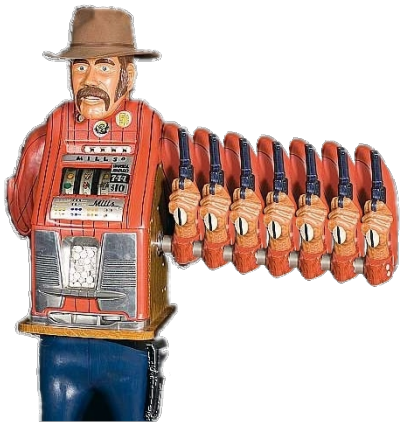
At every decision step $t \in [0, 1, 2, \dots, T]$:

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

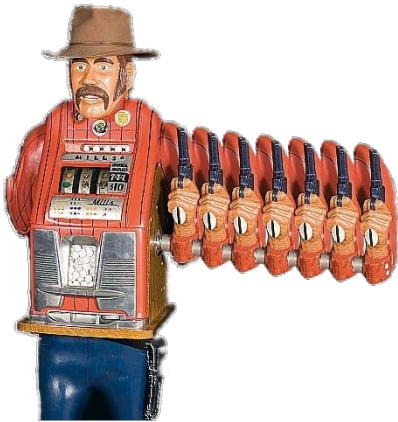
Minimize $\text{Regret}(T) = \sum_{t=0}^T \mu(\theta^*) - \mu(\theta_t)$, where $\theta^* = \arg \max_{\theta \in \Theta} \mu(\theta)$

Problem Formulation



Multi Armed Bandits

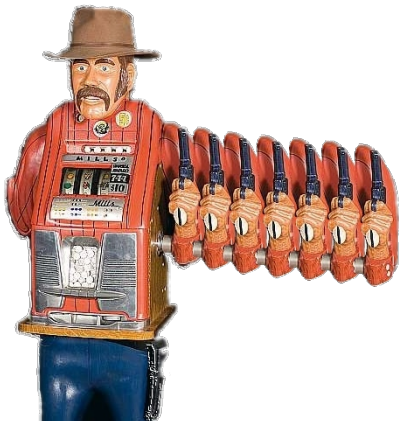
Problem Formulation



Multi Armed Bandits

- Simpler framework;

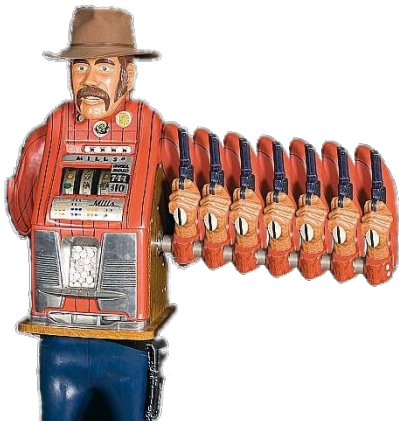
Problem Formulation



Multi Armed Bandits

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

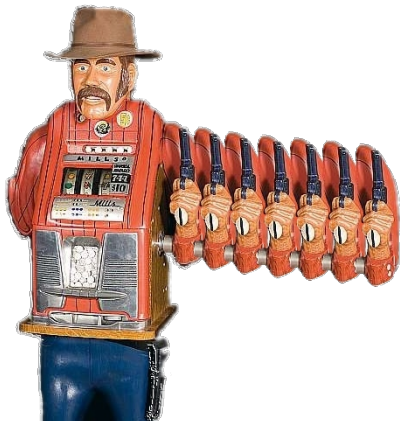
Problem Formulation



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**
5. Experiments
6. Conclusions

Algorithm 1 OPTIMIST

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

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 $\mathcal{R}(\tau_0)$

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**

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 $\theta_t = \arg \max_{\theta \in \Theta} B_t^\epsilon(\theta, \delta_t)$

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**
-

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**
-

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

Discrete Decision Set

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

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

Discrete Decision Set

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

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

Discrete Decision Set

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

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

Discrete Decision Set

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

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

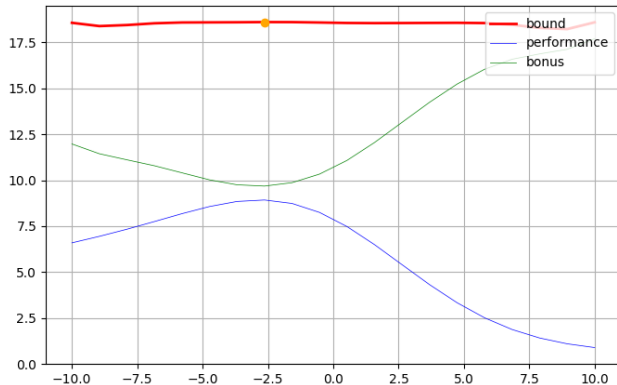
Algorithm 11 OPTIMIST2

- 1: **Input:** initial arm θ_0 , confidence schedule $(\delta_t)_{t=1}^T$, discretization schedule $(\nu_t)_{t=1}^T$, order $\epsilon \in (0, 1]$
 - 2: Draw sample $\nu_0 \sim p_{\theta_0}$ and observe return $f(\tau_0)$
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: Discretize Θ with a uniform grid $\tilde{\Theta}_t$ of ν_t^d points
 - 5: Select arm $\theta_t = \arg \max_{\theta \in \tilde{\Theta}_t} B_t^\epsilon(\theta, \delta_t)$
 - 6: Draw sample $\tau_t \sim p_{\theta_t}$ and observe return $\mathcal{R}(\tau_t)$
 - 7: **end for**
-

Algorithm 12 OPTIMIST2

- 1: **Input:** initial arm θ_0 , confidence schedule $(\delta_t)_{t=1}^T$, discretization schedule $(\nu_t)_{t=1}^T$, order $\epsilon \in (0, 1]$
 - 2: Draw sample $\nu_0 \sim p_{\theta_0}$ and observe return $f(\tau_0)$
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: Discretize Θ with a uniform grid $\tilde{\Theta}_t$ of ν_t^d points
 - 5: Select arm $\theta_t = \arg \max_{\theta \in \tilde{\Theta}_t} B_t^\epsilon(\theta, \delta_t)$
 - 6: Draw sample $\tau_t \sim p_{\theta_t}$ and observe return $\mathcal{R}(\tau_t)$
 - 7: **end for**
-

The Upper Confidence Bound



Theorem (1)

Let \mathcal{X} be a discrete arm set with $|\mathcal{X}| = K \in \mathbb{N}_+$. Under Assumption (??), Algorithm 1 with confidence schedule $\delta_t = \frac{3\delta}{t^2\pi^2K}$ guarantees, with probability at least $1 - \delta$:

$$\text{Regret}(T) \leq \Delta_0 + CT^{\frac{1}{1+\epsilon}} \left[v_\epsilon \left(2 \log T + \log \frac{\pi^2 K}{3\delta} \right) \right]^{\frac{\epsilon}{1+\epsilon}},$$

where $C = (1 + \epsilon) \left(2\sqrt{2} + \frac{5}{3} \right) \|f\|_\infty$, and Δ_0 is the instantaneous regret of the initial arm \mathbf{x}_0 .

This yields a $\mathcal{O}(\sqrt{T \log T})$ regret when $\epsilon = 1$.

Regret Analysis

Theorem (2)

Let \mathcal{X} be a d -dimensional compact arm set with $\mathcal{X} \subseteq [-D, D]^d$. For any $\kappa \geq 2$, under Assumptions (??) and (??), Algorithm 11 with confidence schedule

$\delta_t = \frac{6\delta}{\pi^2 t^2 \left(1 + \lceil t^{1/\kappa} \rceil^d\right)}$ and discretization schedule $\tau_t = \lceil t^{\frac{1}{\kappa}} \rceil$ guarantees, with probability at least $1 - \delta$:

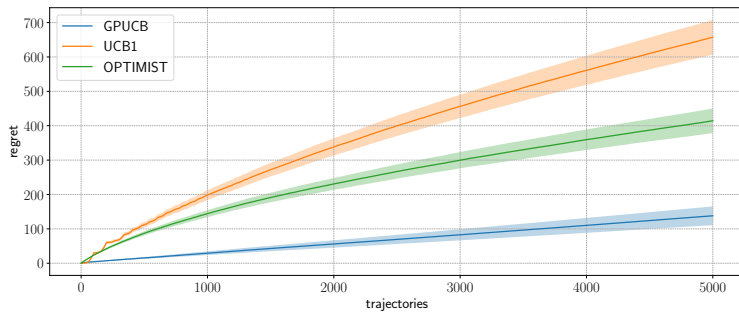
$$\text{Regret}(T) \leq \Delta_0 + C_1 T^{(1-\frac{1}{\kappa})} d + C_2 T^{\frac{1}{1+\epsilon}} \cdot \left[v_\epsilon \left((2 + d/\kappa) \log T + d \log 2 + \log \frac{\pi^2}{3\delta} \right) \right]^{\frac{\epsilon}{1+\epsilon}},$$

where $C_1 = \frac{\kappa}{\kappa - 1} LD$, $C_2 = (1 + \epsilon) \left(2\sqrt{2} + \frac{5}{3} \right) \|f\|_\infty$, and Δ_0 is the instantaneous regret of the initial arm \mathbf{x}_0 .

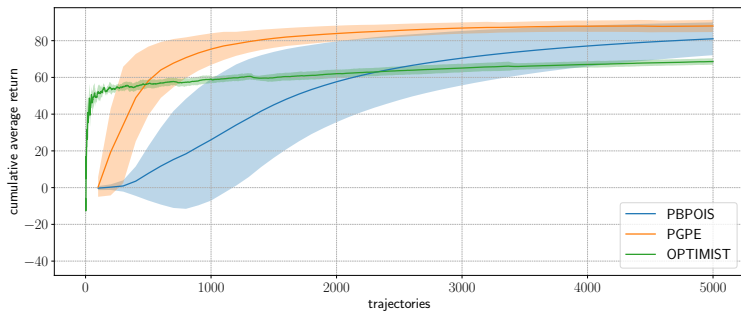
Plan

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

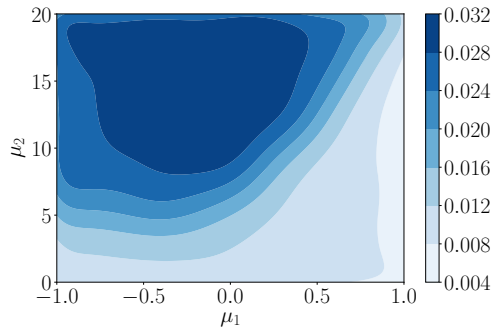
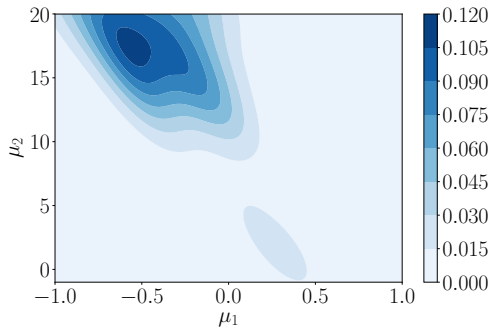
Linear Quadratic Gaussian Regulator



Mountain Car - Performance



Mountain Car - Exploration



Plan

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

Conclusions

Original Contributions

1. Algorithmic contributions: OPTIMIST and OPTIMIST2.

Conclusions

Original Contributions

1. Algorithmic contributions: OPTIMIST and OPTIMIST2.
2. Theoretical contributions:

Conclusions

Original Contributions

1. Algorithmic contributions: OPTIMIST and OPTIMIST2.
2. Theoretical contributions:
 - 2.1 novel problem formalization of Policy Search;

Conclusions

Original Contributions

1. Algorithmic contributions: OPTIMIST and OPTIMIST2.
2. Theoretical contributions:
 - 2.1 novel problem formalization of Policy Search;
 - 2.2 proved sub-linear regret for both algorithms.

Conclusions

Original Contributions

1. Algorithmic contributions: OPTIMIST and OPTIMIST2.
2. Theoretical contributions:
 - 2.1 novel problem formalization of Policy Search;
 - 2.2 proved sub-linear regret for both algorithms.
3. Experimental contributions

Original Contributions

1. Algorithmic contributions: OPTIMIST and OPTIMIST2.
2. Theoretical contributions:
 - 2.1 novel problem formalization of Policy Search;
 - 2.2 proved sub-linear regret for both algorithms.
3. Experimental contributions
 - 3.1 Linear Quadratic Gaussian regulation;

Original Contributions

1. Algorithmic contributions: OPTIMIST and OPTIMIST2.
2. Theoretical contributions:
 - 2.1 novel problem formalization of Policy Search;
 - 2.2 proved sub-linear regret for both algorithms.
3. Experimental contributions
 - 3.1 Linear Quadratic Gaussian regulation;
 - 3.2 Continuous Mountain Car;

Original Contributions

1. Algorithmic contributions: OPTIMIST and OPTIMIST2.
2. Theoretical contributions:
 - 2.1 novel problem formalization of Policy Search;
 - 2.2 proved sub-linear regret for both algorithms.
3. Experimental contributions
 - 3.1 Linear Quadratic Gaussian regulation;
 - 3.2 Continuous Mountain Car;
 - 3.3 **Inverted Pendulum.**

Conclusions

Original Contributions

1. Algorithmic contributions: OPTIMIST and OPTIMIST2.
2. Theoretical contributions:
 - 2.1 novel problem formalization of Policy Search;
 - 2.2 proved sub-linear regret for both algorithms.
3. Experimental contributions
 - 3.1 Linear Quadratic Gaussian regulation;
 - 3.2 Continuous Mountain Car;
 - 3.3 Inverted Pendulum.

Paper submission at **ICML2019** (International Conference on Machine Learning)

Conclusions

Original Contributions

1. Algorithmic contributions: OPTIMIST and OPTIMIST2.
2. Theoretical contributions:
 - 2.1 novel problem formalization of Policy Search;
 - 2.2 proved sub-linear regret for both algorithms.
3. Experimental contributions
 - 3.1 Linear Quadratic Gaussian regulation;
 - 3.2 Continuous Mountain Car;
 - 3.3 Inverted Pendulum.

Paper submission at **ICML2019** (International Conference on Machine Learning)

Future Works

Conclusions

Original Contributions

1. Algorithmic contributions: OPTIMIST and OPTIMIST2.
2. Theoretical contributions:
 - 2.1 novel problem formalization of Policy Search;
 - 2.2 proved sub-linear regret for both algorithms.
3. Experimental contributions
 - 3.1 Linear Quadratic Gaussian regulation;
 - 3.2 Continuous Mountain Car;
 - 3.3 Inverted Pendulum.

Paper submission at **ICML2019** (International Conference on Machine Learning)

Future Works

1. Optimization problem

Conclusions

Original Contributions

1. Algorithmic contributions: OPTIMIST and OPTIMIST2.
2. Theoretical contributions:
 - 2.1 novel problem formalization of Policy Search;
 - 2.2 proved sub-linear regret for both algorithms.
3. Experimental contributions
 - 3.1 Linear Quadratic Gaussian regulation;
 - 3.2 Continuous Mountain Car;
 - 3.3 Inverted Pendulum.






Paper submission at **ICML2019** (International Conference on Machine Learning)

Future Works

1. Optimization problem
2. **Posterior sampling**

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).