



POLITECNICO
MILANO 1863

Exploration in Policy Search by Multiple Importance Sampling

Lorenzo Lupo

Supervisor: Marcello Restelli

Co-supervisors: Matteo Papini, Alberto Maria Metelli

April 16th, 2019

Reinforcement Learning

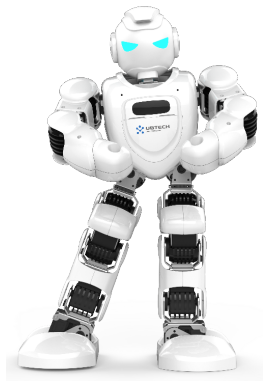
Reinforcement Learning



Reinforcement Learning

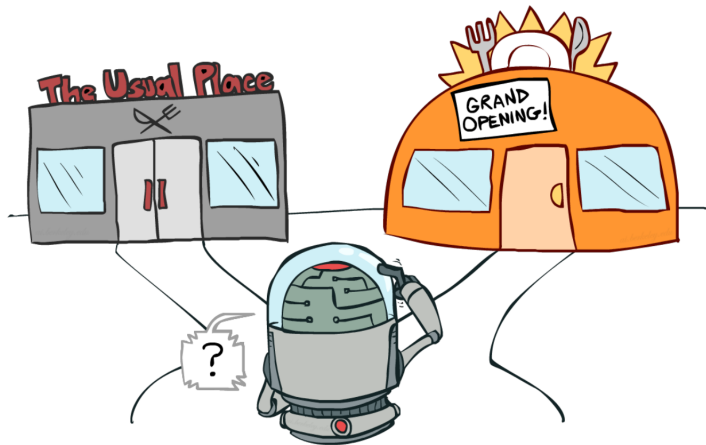


Reinforcement Learning



Exploitation VS Exploration

Exploitation VS Exploration



Plan

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

The Reinforcement Learning Framework

Environment



\mathcal{P}

The Reinforcement Learning Framework

Agent



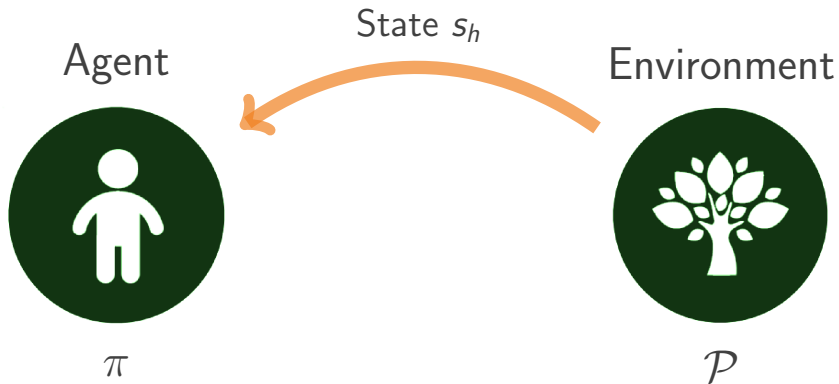
π

Environment

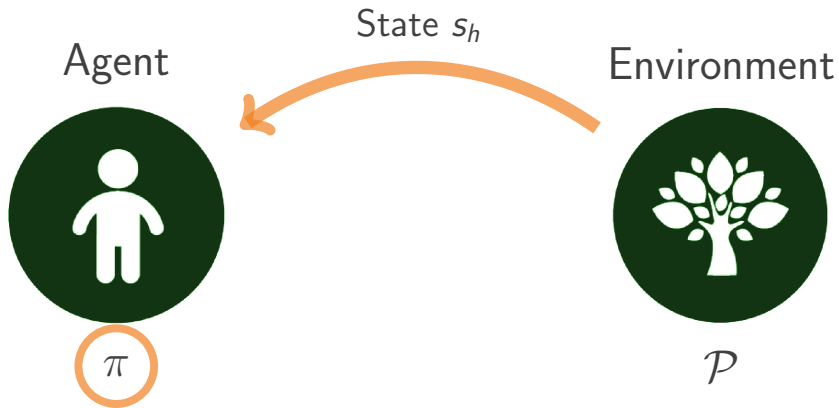


\mathcal{P}

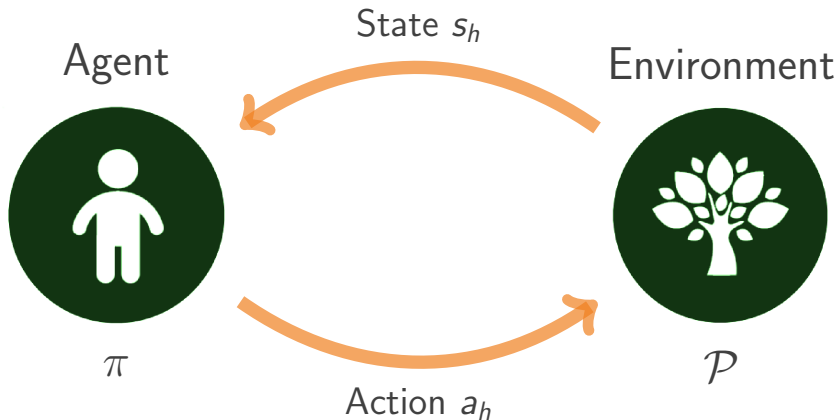
The Reinforcement Learning Framework



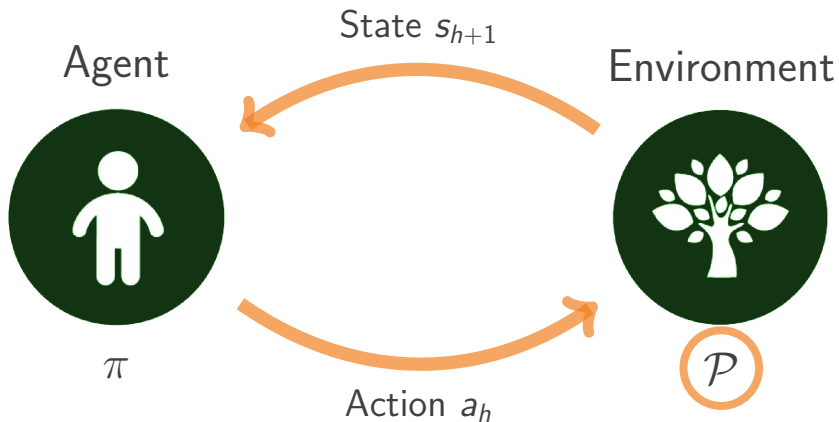
The Reinforcement Learning Framework



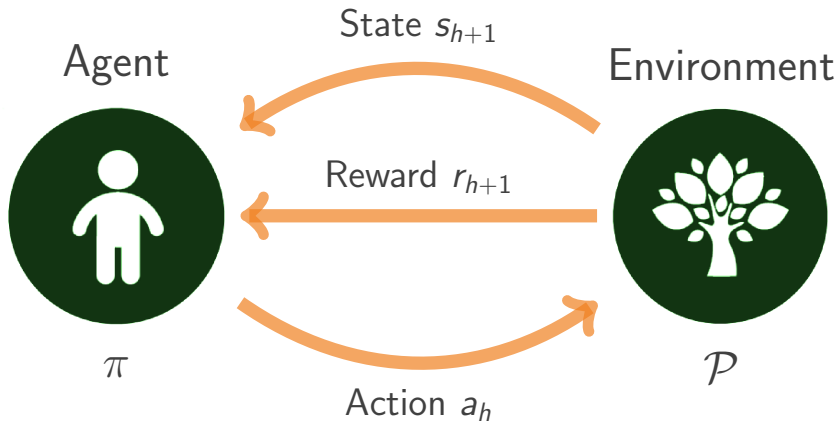
The Reinforcement Learning Framework



The Reinforcement Learning Framework



The Reinforcement Learning Framework



Policy Search Formulation

Parametric policy:

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

Policy Search Formulation

Parametric policy:

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

Return of a trajectory τ :

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

Policy Search Formulation

Parametric policy:

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

Return of a trajectory τ :

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

Performance:

$$\mu(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [\mathcal{R}(\tau)]$$

Policy Search Formulation

Parametric policy:

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

Return of a trajectory τ :

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

Performance:

$$\mu(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [\mathcal{R}(\tau)]$$

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

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

Exploration in Policy Search

Undirected exploration

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

- E.g.1: by adopting stochastic policies [Deisenroth et al., 2013].

Exploration in Policy Search

Undirected exploration

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

- E.g.1: by adopting stochastic policies [Deisenroth et al., 2013].
- E.g.2: by augmenting rewards with the entropy of the policy [Haarnoja et al., 2018]:

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

Exploration in Policy Search

Undirected exploration

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

- E.g.1: by adopting stochastic policies [Deisenroth et al., 2013].
- E.g.2: by augmenting rewards with the entropy of the policy [Haarnoja et al., 2018]:

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

Directed exploration

Leverage on the knowledge acquired during learning.

Exploration in Policy Search

Undirected exploration

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

- E.g.1: by adopting stochastic policies [Deisenroth et al., 2013].
- E.g.2: by augmenting rewards with the entropy of the policy [Haarnoja et al., 2018]:

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

Directed exploration

Leverage on the knowledge acquired during learning.

- E.g.: Count-based techniques [Bellemare et al., 2016].

Plan

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

Problem Formalization

Decision Set

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

Problem Formalization

Decision 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

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

Procedure

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

1. **Select** a parametrization $\theta_t \in \Theta$;

Problem Formalization

Decision Set

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

Procedure

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

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

Problem Formalization

Decision Set

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

Procedure

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

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

Problem Formalization

Decision Set

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

Procedure

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

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

Problem Formalization

Decision Set

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

Procedure

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

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

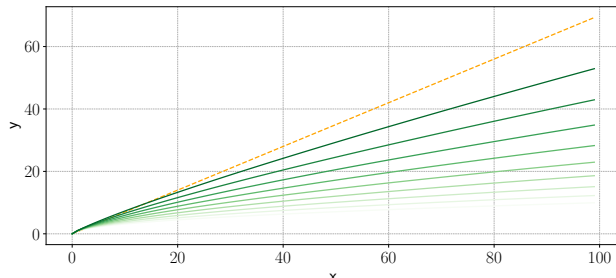
Goal

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

Desideratum

Sub-linear Regret

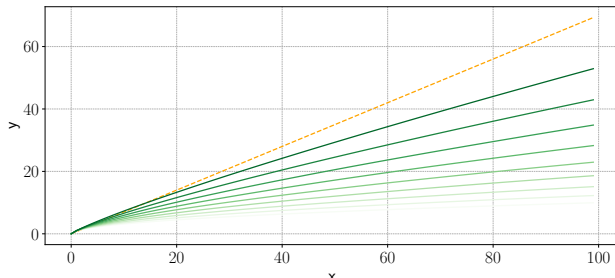
1. for any linear function $f(T)$, for sufficiently large input T , $\text{Regret}(T)$ grows slower than $f(T)$;



Desideratum

Sub-linear Regret

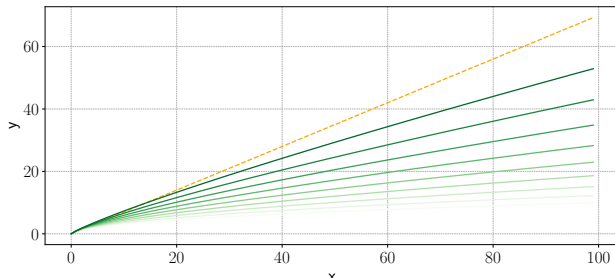
1. for any linear function $f(T)$, for sufficiently large input T , $\text{Regret}(T)$ grows slower than $f(T)$;
2. Alternatively, $\lim_{T \rightarrow \infty} \text{Regret}(T)/T = 0$.



Desideratum

Sub-linear Regret

1. for any linear function $f(T)$, for sufficiently large input T , $\text{Regret}(T)$ grows slower than $f(T)$;
2. Alternatively, $\lim_{T \rightarrow \infty} \text{Regret}(T)/T = 0$.
3. **Meaning:** after a certain number of iterations, the policy keeps improving.



Plan

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

Optimistic Policy Optimization via Multiple Importance Sampling with Truncation

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

Optimistic Policy Optimization via Multiple Importance Sampling with Truncation

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

$$\mu(\theta) \leq B_t(\theta, \delta_t)$$

Optimistic Policy Optimization via Multiple Importance Sampling with Truncation

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

$$\mu(\theta) \leq B_t(\theta, \delta_t)$$

$$B_t(\theta, \delta_t) = \check{\mu}_t^{MIS}(\theta) + C_b \sqrt{\frac{d_2(\pi_{\theta_t} \parallel \Phi_t) \log \frac{1}{\delta_t}}{t}}$$

Optimistic Policy Optimization via Multiple Importance Sampling with Truncation

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

$$\mu(\theta) \leq B_t(\theta, \delta_t)$$

Truncated Multiple Importance Sampling Estimator

$$B_t(\theta, \delta_t) = \boxed{\check{\mu}_t^{MIS}(\theta)} + C_b \sqrt{\frac{d_2(\pi_{\theta_t} \parallel \Phi_t) \log \frac{1}{\delta_t}}{t}}$$

Optimistic Policy Optimization via Multiple Importance Sampling with Truncation

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

$$\mu(\theta) \leq B_t(\theta, \delta_t)$$

Exploration Bonus

$$B_t(\theta, \delta_t) = \check{\mu}_t^{MIS}(\theta) + C_b \sqrt{\frac{d_2(\pi_{\theta_t} \parallel \Phi_t) \log \frac{1}{\delta_t}}{t}}$$

Optimistic Policy Optimization via Multiple Importance Sampling with Truncation

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

$$\mu(\theta) \leq B_t(\theta, \delta_t)$$

Exploration Bonus

$$B_t(\theta, \delta_t) = \check{\mu}_t^{MIS}(\theta) + C_b \sqrt{\frac{d_2(\pi_{\theta_t} \parallel \Phi_t) \log \frac{1}{\delta_t}}{t}}$$

Regret Analysis - Discrete Arms Set

Theorem 1 - Discrete Parameter Space:

$$\text{Regret}(T) \leq \Delta_0 + C_1 \sqrt{T \left[v_1 \left(2 \log T + \log \frac{2\pi^2}{3\delta} \right) \right]}$$

Regret Analysis - Discrete Arms Set

Theorem 1 - Discrete Parameter Space:

$$\text{Regret}(T) \leq \Delta_0 + C_1 \sqrt{T \left[v_1 \left(2 \log T + \log \frac{2\pi^2}{3\delta} \right) \right]} = \mathcal{O}(\sqrt{T \log T})$$

Regret Analysis - Discrete Arms Set

Theorem 1 - Discrete Parameter Space:

$$\text{Regret}(T) \leq \Delta_0 + C_1 \sqrt{T \left[v_1 \left(2 \log T + \log \frac{2\pi^2}{3\delta} \right) \right]} = \mathcal{O}(\sqrt{T \log T})$$

Theorem 2 - Compact Parameter Space:

$$\text{Regret}(T) \leq C_2 + C_3 \sqrt{T \left[v_1 \left(2(d+1) \log T + d \log d + \log \frac{\pi^2}{3\delta} \right) \right]}$$

Regret Analysis - Discrete Arms Set

Theorem 1 - Discrete Parameter Space:

$$\text{Regret}(T) \leq \Delta_0 + C_1 \sqrt{T \left[v_1 \left(2 \log T + \log \frac{2\pi^2}{3\delta} \right) \right]} = \mathcal{O}(\sqrt{T \log T})$$

Theorem 2 - Compact Parameter Space:

$$\text{Regret}(T) \leq C_2 + C_3 \sqrt{T \left[v_1 \left(2(d+1) \log T + d \log d + \log \frac{\pi^2}{3\delta} \right) \right]} = \mathcal{O}(\sqrt{dT \log T})$$

Regret Analysis - Discrete Arms Set

Theorem 1 - Discrete Parameter Space:

$$\text{Regret}(T) \leq \Delta_0 + C_1 \sqrt{T \left[v_1 \left(2 \log T + \log \frac{2\pi^2}{3\delta} \right) \right]} = \mathcal{O}(\sqrt{T \log T})$$

Theorem 2 - Compact Parameter Space:

$$\text{Regret}(T) \leq C_2 + C_3 \sqrt{T \left[v_1 \left(2(d+1) \log T + d \log d + \log \frac{\pi^2}{3\delta} \right) \right]} = \mathcal{O}(\sqrt{dT \log T})$$

Theorem 3 - Discretized Parameter Space:

$$\text{Reg}(T) \leq \Delta_0 + C_4 \sqrt{T} d + C_5 \sqrt{T \left[v_1 \left((2 + d/2) \log T + d \log 2 + \log \frac{\pi^2}{3\delta} \right) \right]}$$

Regret Analysis - Discrete Arms Set

Theorem 1 - Discrete Parameter Space:

$$\text{Regret}(T) \leq \Delta_0 + C_1 \sqrt{T \left[v_1 \left(2 \log T + \log \frac{2\pi^2}{3\delta} \right) \right]} = \mathcal{O}(\sqrt{T \log T})$$

Theorem 2 - Compact Parameter Space:

$$\text{Regret}(T) \leq C_2 + C_3 \sqrt{T \left[v_1 \left(2(d+1) \log T + d \log d + \log \frac{\pi^2}{3\delta} \right) \right]} = \mathcal{O}(\sqrt{dT \log T})$$

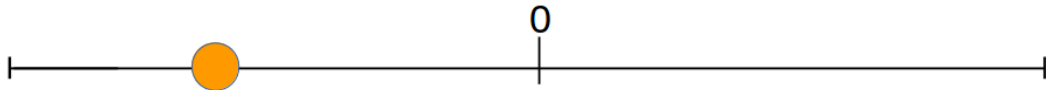
Theorem 3 - Discretized Parameter Space:

$$\text{Reg}(T) \leq \Delta_0 + C_4 \sqrt{T} d + C_5 \sqrt{T \left[v_1 \left((2 + d/2) \log T + d \log 2 + \log \frac{\pi^2}{3\delta} \right) \right]} = \mathcal{O}(\sqrt{dT \log T})$$

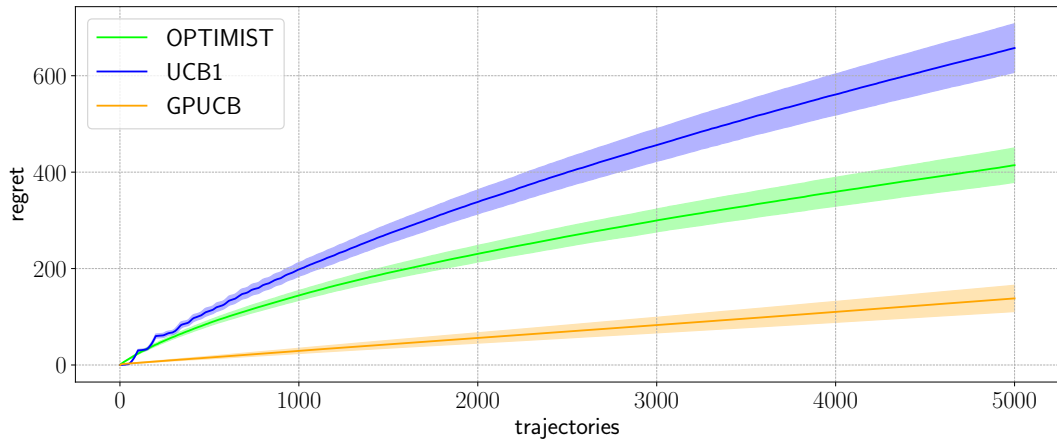
Plan

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

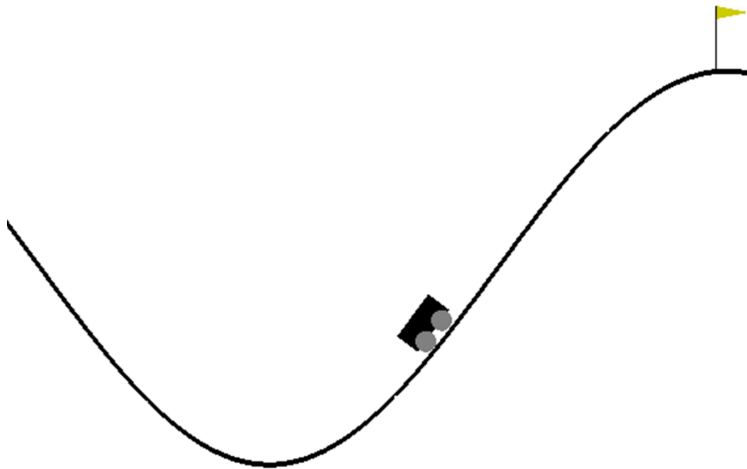
Linear Quadratic Gaussian Regulator



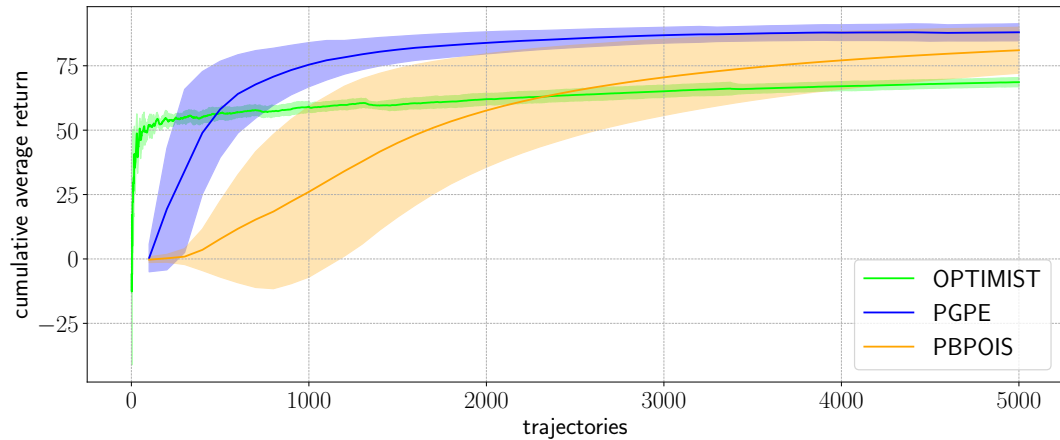
Linear Quadratic Gaussian Regulator - Regret



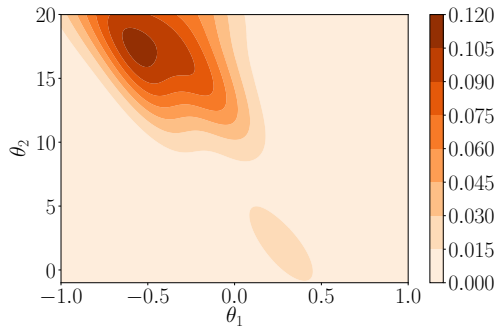
Mountain Car



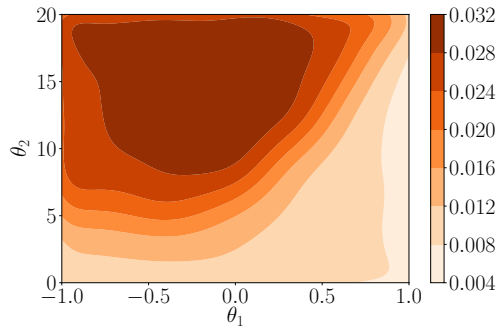
Mountain Car - Performance



Mountain Car - Parameter Space Exploration



PGPE



OPTIMIST

Plan

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

Original Contributions

1. Algorithmic contributions: OPTIMIST

Original Contributions

1. Algorithmic contributions: OPTIMIST
2. Theoretical contributions:

Original Contributions

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

Original Contributions

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

Original Contributions

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

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

Conclusions

Original Contributions

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

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

Future Works

Conclusions

Original Contributions

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

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

Future Works

1. Optimization problem;

Conclusions

Original Contributions

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




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

Future Works

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

Thank you for your attention!

References I

-  Bellemare, M., Srinivasan, S., Ostrovski, G., Schaul, T., Saxton, D., and Munos, R. (2016).
Unifying count-based exploration and intrinsic motivation.
In Advances in Neural Information Processing Systems, pages 1471–1479.
-  Bubeck, S., Cesa-Bianchi, N., and Lugosi, G. (2013).
Bandits with heavy tail.
IEEE Transactions on Information Theory, 59(11):7711–7717.
-  Deisenroth, M. P., Neumann, G., Peters, J., et al. (2013).
A survey on policy search for robotics.
Foundations and Trends® in Robotics, 2(1–2):1–142.

References II

-  Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. (2018).
Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor.
In Proceedings of the 35th International Conference on Machine Learning, pages 1856–1865.
-  Lattimore, T. and Szepesvári, C. (2019).
Bandit Algorithms.
Cambridge University Press (preprint).
-  Srinivas, N., Krause, A., Kakade, S., and Seeger, M. (2010).
Gaussian process optimization in the bandit setting: No regret and experimental design.

References III

In Fürnkranz, J. and Joachims, T., editors, *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 1015–1022, Haifa, Israel. Omnipress.



Sutton, R. S. and Barto, A. G. (2018).
Reinforcement learning: An introduction.
MIT press.