



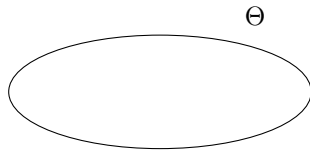
POLITECNICO
MILANO 1863

Optimistic Policy Optimization via Multiple Importance Sampling

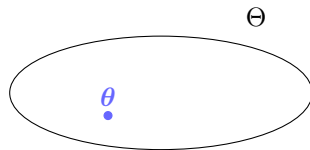
Matteo Papini Alberto Maria Metelli
Lorenzo Lupo Marcello Restelli

11th June 2019

Thirty-sixth International Conference on Machine Learning, Long Beach, CA, USA

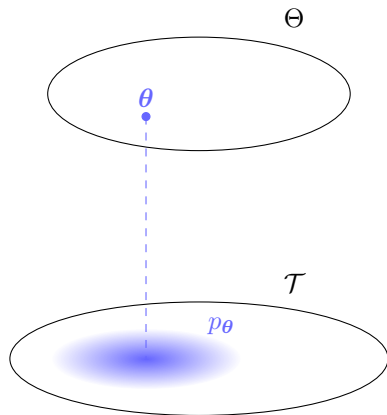


- **Parameter space** $\Theta \subseteq \mathbb{R}^d$
- A parametric **policy** for each $\theta \in \Theta$
- Each inducing a distribution p_θ over **trajectories**
- A **return** $R(\tau)$ for every trajectory τ
- **Goal:** $\max_{\theta \in \Theta} J(\theta) = \mathbb{E}_{\tau \sim p_\theta} [R(\tau)]$
- Iterative optimization (e.g., gradient ascent)

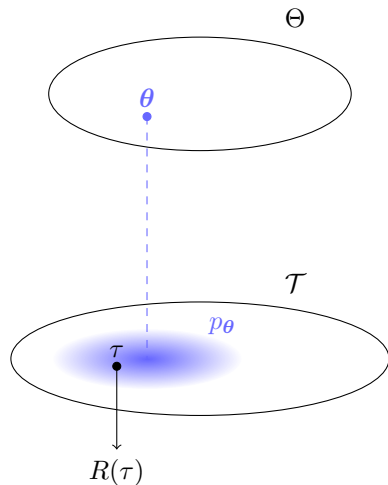


- **Parameter space** $\Theta \subseteq \mathbb{R}^d$
- A parametric **policy** for each $\theta \in \Theta$
- Each inducing a distribution p_θ over **trajectories**
- A **return** $R(\tau)$ for every trajectory τ
- **Goal:** $\max_{\theta \in \Theta} J(\theta) = \mathbb{E}_{\tau \sim p_\theta} [R(\tau)]$
- Iterative optimization (e.g., gradient ascent)

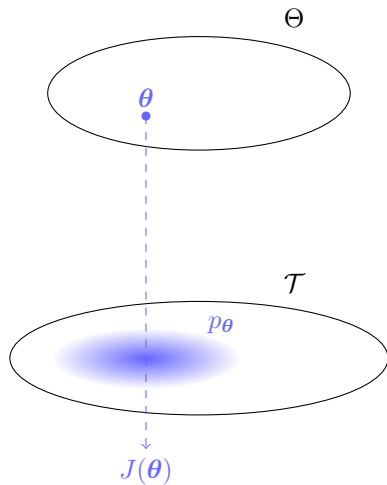
- **Parameter space** $\Theta \subseteq \mathbb{R}^d$
- A parametric **policy** for each $\theta \in \Theta$
- Each inducing a distribution p_θ over **trajectories**
- A **return** $R(\tau)$ for every trajectory τ
- **Goal:** $\max_{\theta \in \Theta} J(\theta) = \mathbb{E}_{\tau \sim p_\theta} [R(\tau)]$
- Iterative optimization (e.g., gradient ascent)



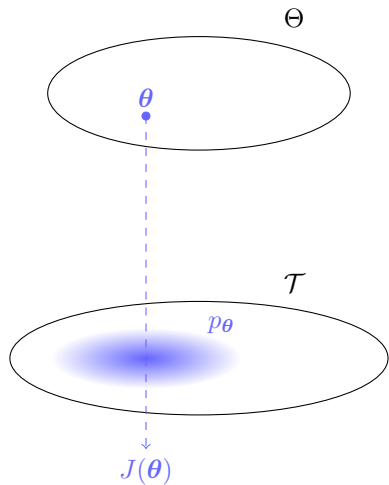
- **Parameter space** $\Theta \subseteq \mathbb{R}^d$
- A parametric **policy** for each $\theta \in \Theta$
- Each inducing a distribution p_θ over **trajectories**
- A **return** $R(\tau)$ for every trajectory τ
- **Goal:** $\max_{\theta \in \Theta} J(\theta) = \mathbb{E}_{\tau \sim p_\theta} [R(\tau)]$
- Iterative optimization (e.g., gradient ascent)



- **Parameter space** $\Theta \subseteq \mathbb{R}^d$
- A parametric **policy** for each $\theta \in \Theta$
- Each inducing a distribution p_θ over **trajectories**
- A **return** $R(\tau)$ for every trajectory τ
- **Goal:** $\max_{\theta \in \Theta} J(\theta) = \mathbb{E}_{\tau \sim p_\theta} [R(\tau)]$
- Iterative optimization (e.g., gradient ascent)



- **Parameter space** $\Theta \subseteq \mathbb{R}^d$
- A parametric **policy** for each $\theta \in \Theta$
- Each inducing a distribution p_θ over **trajectories**
- A **return** $R(\tau)$ for every trajectory τ
- **Goal:** $\max_{\theta \in \Theta} J(\theta) = \mathbb{E}_{\tau \sim p_\theta} [R(\tau)]$
- Iterative optimization (e.g., gradient ascent)



- **Exploration-exploitation** trade-off
- The underlying Markov process is often **continuous**
- **Undirected** exploration: entropy bonus [3]
- **Directed** exploration: pseudo-counts [1]

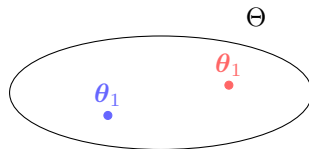
- **Exploration-exploitation** trade-off
- The underlying Markov process is often **continuous**
- **Undirected** exploration: entropy bonus [3]
- **Directed** exploration: pseudo-counts [1]

- **Exploration-exploitation** trade-off
- The underlying Markov process is often **continuous**
- **Undirected** exploration: entropy bonus [3]
- **Directed** exploration: pseudo-counts [1]

- **Exploration-exploitation** trade-off
- The underlying Markov process is often **continuous**
- **Undirected** exploration: entropy bonus [3]
- **Directed** exploration: pseudo-counts [1]

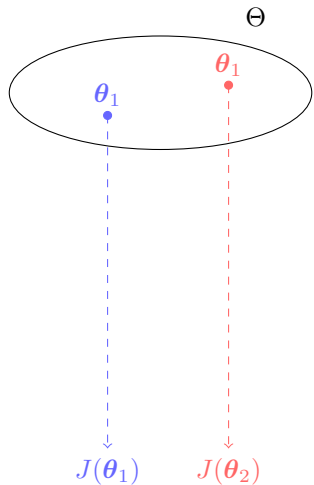
- **Exploration-exploitation** trade-off
- The underlying Markov process is often **continuous**
- **Undirected** exploration: entropy bonus [3]
- **Directed** exploration: pseudo-counts [1]

Lack of theoretical guarantees

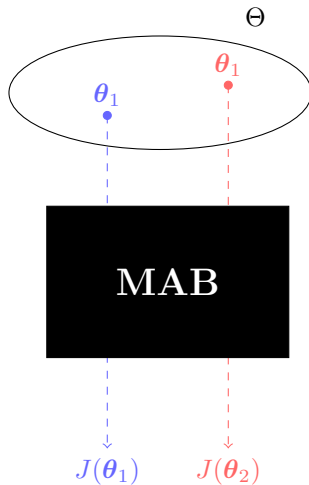


- **Arms:** parameters θ
- **Payoff:** expected return $J(\theta)$
- **Continuous MAB** [4]: we *need* structure
- **Arm correlation** [6] through trajectory distributions

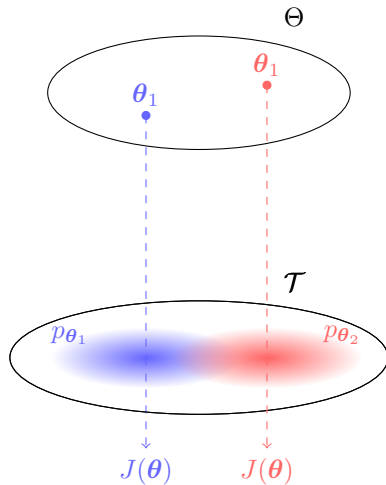
- **Arms:** parameters θ
- **Payoff:** expected return $J(\theta)$
- **Continuous MAB** [4]
- **Arm correlation** [6] through trajectory distributions



- **Arms:** parameters θ
- **Payoff:** expected return $J(\theta)$
- **Continuous MAB** [4]
- **Arm correlation** [6] through trajectory distributions



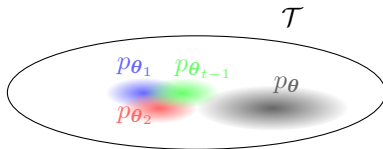
- **Arms:** parameters θ
- **Payoff:** expected return $J(\theta)$
- **Continuous MAB** [4]
- **Arm correlation** [6] through trajectory distributions



- A **UCB-like** index [5]:

$$B_t(\boldsymbol{\theta}) = \underbrace{\check{J}_t(\boldsymbol{\theta})}_{\text{ESTIMATE}}$$

a **truncated multiple**
importance sampling estimator [8, 2]

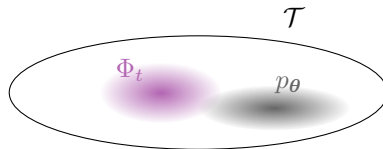


- A **UCB-like** index [5]:

$$B_t(\theta) = \underbrace{\check{J}_t(\theta)}_{\text{ESTIMATE}} + \underbrace{C \sqrt{\frac{d_2(p_\theta \| \Phi_t) \log \frac{1}{\delta_t}}{t}}}_{\text{EXPLORATION BONUS:}}$$

a **truncated multiple** importance sampling estimator [8, 2]

distributional distance from previous solutions



- A **UCB-like** index [5]:

$$B_t(\boldsymbol{\theta}) = \underbrace{\check{J}_t(\boldsymbol{\theta})}_{\text{ESTIMATE}} + \underbrace{C \sqrt{\frac{d_2(p_{\boldsymbol{\theta}} \parallel \Phi_t) \log \frac{1}{\delta_t}}{t}}}_{\text{EXPLORATION BONUS:}}$$

a **truncated multiple**
importance sampling estimator [8, 2]

distributional distance
from previous solutions

- Select $\boldsymbol{\theta}_t = \arg \max_{\boldsymbol{\theta} \in \Theta} B_t(\boldsymbol{\theta})$

- $Regret(T) = \sum_{t=0}^T J(\boldsymbol{\theta}^*) - J(\boldsymbol{\theta}_t)$

- **Compact**, d -dimensional parameter space Θ

- Under **mild assumptions** on the policy class, with high probability:

$$Regret(T) = \tilde{\mathcal{O}}\left(\sqrt{dT}\right)$$

- $Regret(T) = \sum_{t=0}^T J(\boldsymbol{\theta}^*) - J(\boldsymbol{\theta}_t)$
- **Compact**, d -dimensional parameter space Θ

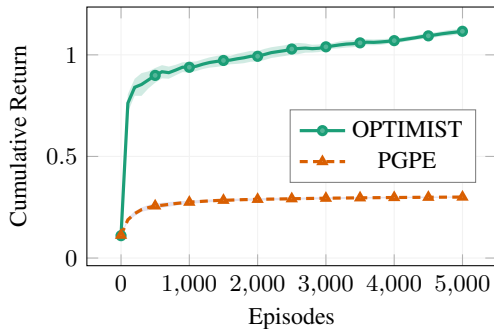
- Under **mild assumptions** on the policy class, with high probability:

$$Regret(T) = \tilde{\mathcal{O}}\left(\sqrt{dT}\right)$$

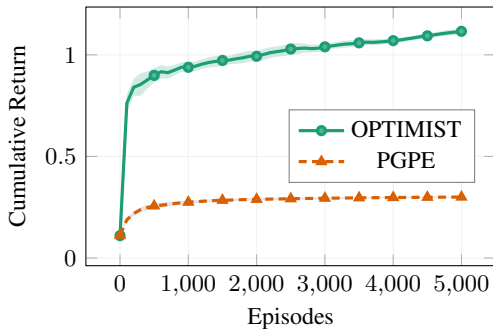
- $Regret(T) = \sum_{t=0}^T J(\boldsymbol{\theta}^*) - J(\boldsymbol{\theta}_t)$
- **Compact**, d -dimensional parameter space Θ
- Under **mild assumptions** on the policy class, with high probability:

$$Regret(T) = \tilde{\mathcal{O}}\left(\sqrt{dT}\right)$$

River Swim



River Swim



Caveats

- Easy implementation only for parameter-based exploration [7]
- Difficult optimization \Rightarrow discretization
- ...

Thank You for Your Attention!

Poster **#103**

Code: `github.com/WolfLo/optimist`

Contact: `matteo.papini@polimi.it`

Web page: `t3p.github.io/icml19`



- [1] 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.
- [2] Bubeck, S., Cesa-Bianchi, N., and Lugosi, G. (2013). Bandits with heavy tail. *IEEE Transactions on Information Theory*, 59(11):7711–7717.
- [3] 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, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, pages 1856–1865.
- [4] Kleinberg, R., Slivkins, A., and Upfal, E. (2013). Bandits and experts in metric spaces. *arXiv preprint arXiv:1312.1277*.
- [5] Lai, T. L. and Robbins, H. (1985). Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22.
- [6] Pandey, S., Chakrabarti, D., and Agarwal, D. (2007). Multi-armed bandit problems with dependent arms. In *Proceedings of the 24th international conference on Machine learning*, pages 721–728. ACM.
- [7] Sehnke, F., Osendorfer, C., Rückstieß, T., Graves, A., Peters, J., and Schmidhuber, J. (2008). Policy gradients with parameter-based exploration for control. In *International Conference on Artificial Neural Networks*, pages 387–396. Springer.
- [8] Veach, E. and Guibas, L. J. (1995). Optimally combining sampling techniques for Monte Carlo rendering. In *Proceedings of the 22nd annual conference on Computer graphics and interactive techniques - SIGGRAPH '95*, pages 419–428. ACM Press.