



POLITECNICO
MILANO 1863

OPTIMISTIC POLICY OPTIMIZATION VIA MULTIPLE IMPORTANCE SAMPLING

MATTEO PAPINI, ALBERTO M. METELLI, LORENZO LUPO AND MARCELLO RESTELLI
{matteo.papini, albertomaria.metelli, marcello.restelli}@polimi.it, lorenzo.lupo@mail.polimi.it



MOTIVATION AND IDEA

Problem:

- Policy Optimization (PO) methods **neglect exploration**
- Existing exploration strategies are **undirected**
- Lack of **provably efficient** solutions

Idea:

- Frame PO as a **Multi-Armed Bandit (MAB)** over parameter space *with **arm correlation***
- Apply **Optimism in Face of Uncertainty (OFU)**

POLICY OPTIMIZATION

Vanilla action-based PO (Peters and Schaal, 2008)

- **Continuous** MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \mu \rangle$
- **Trajectory** $\tau = s_0, a_0, r_1, s_1, \dots, r_H \in \mathcal{T}$
- Return $\mathcal{R}(\tau) = \sum_{h=0}^{H-1} \gamma^h r_{h+1}$
- **Parametric** policy $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$ with $\theta \in \Theta$
- Induced **trajectory distribution** $p_\theta(\tau)$
- **Performance** $J(\theta) = \mathbb{E}_{\tau \sim p_\theta}[\mathcal{R}(\tau)]$
- Find $\theta^* = \arg \max_{\theta \in \Theta} J(\theta)$

Parameter-based PO (Sehnke et al., 2008):

- **Hyperpolicy** $\nu_\xi(\theta)$ with $\xi \in \Xi$ (e.g., Gaussian)
- Find $\xi^* = \arg \max_{\xi \in \Xi} \mathbb{E}_{\theta \sim \nu_\xi} [J(\theta)]$

POLICY OPTIMIZATION AS CORRELATED MAB

	Correlated MAB	PO	PB-PO
Arm	$x \in \mathcal{X}$	$\theta \in \Theta$	$\xi \in \Xi$
Outcome	$z \in \mathcal{Z}$	$\tau \in \mathcal{T}$	$\theta \in \Theta$
Induced distribution	$p_x(z)$	$p_\theta(\tau)$	$\nu_\xi(\theta)$
Payoff	$f(z)$	$\mathcal{R}(\tau)$	$J(\theta)$

MULTIPLE IMPORTANCE SAMPLING

- Samples from several **behavioral** distributions:
 $z_0 \sim q_0, z_1 \sim q_1, \dots, z_K \sim q_K$
- Estimate $\mu := \mathbb{E}_{z \sim p} [f(z)]$ under **target** distribution p
- **Balance Heuristic (BH)** (Veach and Guibas, 1995):

$$\hat{\mu}_{\text{BH}} := \frac{1}{K} \sum_{k=1}^K \underbrace{\frac{p(z_k)}{\Phi(z_k)}}_{\text{Importance Weight (IW)}} f(z_k), \quad \Phi(z) = \underbrace{\frac{1}{K} \sum_{k=1}^K q_k(z)}_{\text{mixture}}$$

- **Unbiased**, but possibly **high-variance**:

$$\mathbb{V}\text{ar} [\hat{\mu}_{\text{BH}}] \leq \|f\|_\infty^2 \frac{d_2(P\|\Phi)}{K} \leq \|f\|_\infty^2 \frac{1}{\sum_{k=1}^K \frac{1}{d_2(p\|q_k)}}$$
$$d_2(p\|q) := \int_{\mathcal{Z}} \left(\frac{p(z)}{q(z)} \right)^2 dz \quad (\text{Rényi divergence})$$

OPTIMIST ALGORITHM

ROBUST ESTIMATOR

- Importance Sampling estimators are **heavy-tailed** (Metelli et al., 2018)
- This prevents the formation of *exponential* **Upper Confidence Bounds (UCB)**
- Robust estimation via **adaptive truncation** (Bubeck et al., 2012):

$$\check{\mu}_{\text{BH}} := \frac{1}{K} \sum_{k=1}^K \min \left\{ \underbrace{\sqrt{\frac{K d_2(p\|\Phi)}{\log \frac{1}{\delta}}}}_{\text{truncation}}, \underbrace{\frac{p(z_k)}{\Phi(z_k)}}_{\text{IW}} \right\} f(z_k)$$

- Thanks to truncation, with probability at least $1 - 2\delta$:

$$|\check{\mu}_{\text{BH}} - \mu| \leq \|f\|_\infty \left(\sqrt{2} + \frac{4}{3} \right) \sqrt{\frac{d_2(p\|\Phi) \log \frac{1}{\delta}}{K}}$$

IMPLEMENTATION

EXPERIMENTS

REFERENCES

S. Bubeck, N. Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.
A. M. Metelli, M. Papini, F. Faccio, and M. Restelli. Policy optimization via importance sampling. In *Advances in Neural Information Processing Systems*, pages 5447–5459, 2018.
J. Peters and S. Schaal. Reinforcement learning of motor skills with policy gradients. *Neural networks*, 21(4):682–697, 2008.
F. Sehnke, C. Osendorfer, T. Rückstieß, A. Graves, J. Peters, and J. Schmidhuber. Policy gradients with parameter-based exploration for control. In *International Conference on Artificial Neural Networks*, pages 387–396. Springer, 2008.
E. Veach and L. J. Guibas. 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, 1995.

REGRET ANALYSIS