



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 a lot of structure*

POLICY OPTIMIZATION

Vanilla:

- **Continuous** MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \mu \rangle$
- **Trajectory** $\tau = s_0, a_0, r_1, s_1, \dots, r_H$
- 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 exploration (Sehnke et al., 2008):

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

MULTIPLE IMPORTANCE SAMPLING

- Samples from many **behavioral distributions**:
 $z_0 \sim q_0, z_1 \sim q_1, \dots, z_K \sim q_K$
- **Target distribution** p
- Estimate $\mu = \mathbb{E}_{z \sim p} [f(z)]$ from available samples

- **Mixture** of behaviorals: $\phi(z) = \frac{1}{K} \sum_{k=1}^K q_k(z)$

- **Balance Heuristic (BH)** (Veach and Guibas, 1995):

$$\hat{\mu}_{\text{BH}} := \frac{1}{K} \sum_{k=1}^K \frac{p(z_k)}{\phi(z_k)} f(z_k) \quad (\text{unbiased})$$

- Possibly **high variance**:

$$\text{Var} [\hat{\mu}_{\text{BH}}] \leq \|f\|_\infty^2 \frac{d_2(P \parallel \Phi)}{K} \leq \frac{\|f\|_\infty^2}{\sum_{k=1}^K \frac{1}{d_2(p \parallel q_k)}}$$

$$d_2(p \parallel q) := \int_{\mathcal{Z}} \left(\frac{p(z)}{q(z)} \right)^2 \text{d}z \quad (\text{exp. Renyi divergence})$$

ROBUST MULTIPLE IMPORTANCE SAMPLING ESTIMATION

ALGORITHM

REGRET ANALYSIS

IMPLEMENTATION

EXPERIMENTS

REFERENCES

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.