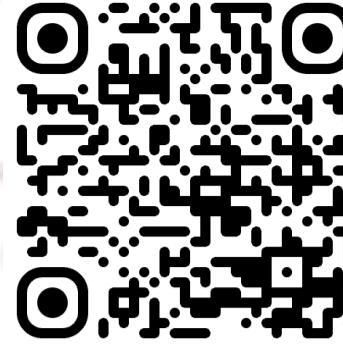


CONVERGENCE ANALYSIS OF POLICY GRADIENT METHODS WITH DYNAMIC STOCHASTICITY

A. MONTENEGRO, M. MUSSI, M. PAPINI, AND A. M. METELLI
{alessandro.montenegro, marco.mussi, matteo.papini, albertomaria.metelli}@polimi.it



MOTIVATION AND CONTRIBUTION

WHY POLICY GRADIENTS?

They successfully tackle real-world **continuous control problems**.

WHY DETERMINISTIC POLICIES?

Real-life agents cannot accept stochastic policies, they do not meet **reliability**, **safety**, and **traceability** standards.

WHY DYNAMIC STOCHASTICITY?

Theorists consider it **fixed** to show convergence, but to deploy deterministic policies it has to be set small. Practitioners learn it, not caring of convergence.

PES: DECREASE THE STOCHASTICITY.

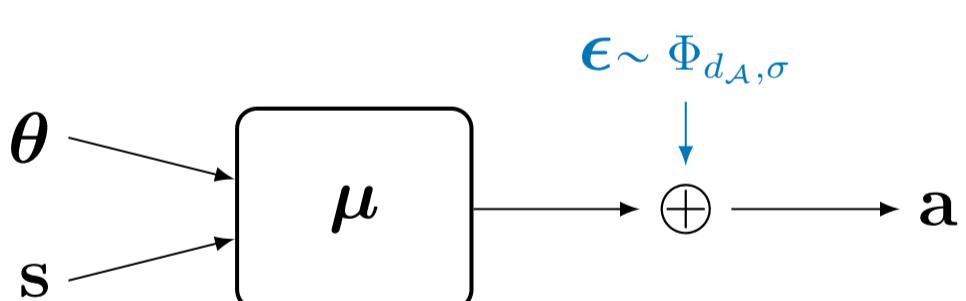
Convergence to the optimal **deterministic** policy in $\tilde{\mathcal{O}}(\epsilon^{-5})$, as if keeping a small fixed stochasticity!

SL-PG: LEARN THE STOCHASTICITY.

Learn a stochasticity mapping in $[\sigma_{\min}, \sigma_{\max}]$. Convergence to the optimal **stochastic** policy in $\tilde{\mathcal{O}}(\epsilon^{-3}\sigma_{\min}^{-2})$, but stronger assumptions are required!

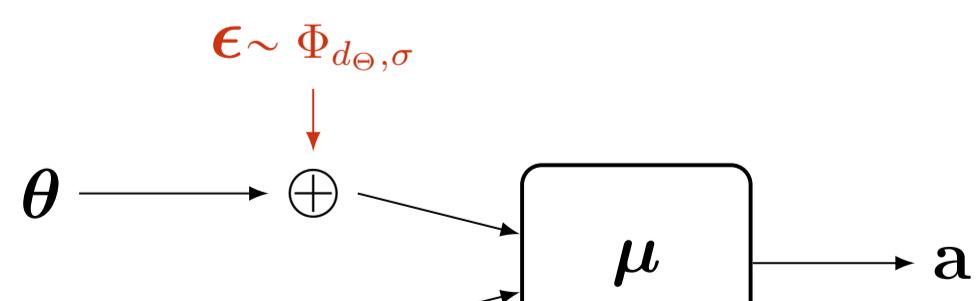
AB AND PB EXPLORATION

ACTION-BASED EXPLORATION



$$J_A(\theta, \sigma) = \mathbb{E}_{\tau \sim p_A(\cdot | \theta, \sigma)} [R(\tau)]$$

PARAMETER-BASED EXPLORATION



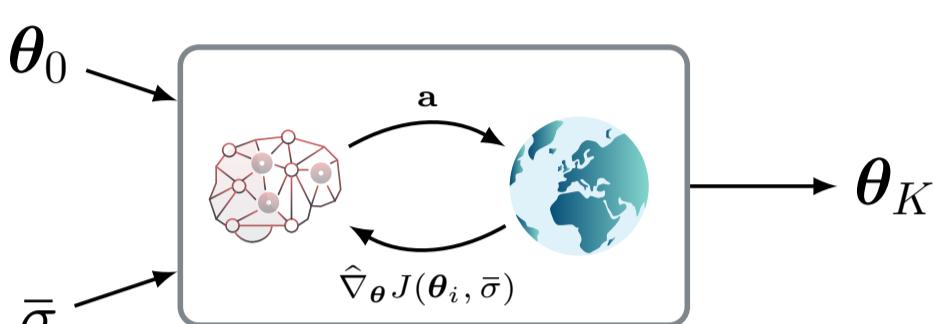
$$J_P(\theta, \sigma) = \mathbb{E}_{\theta \sim \nu_{P, \sigma}} \left[\mathbb{E}_{\tau \sim p_D(\cdot | \theta)} [R(\tau)] \right]$$

CONVERGENCE: FIXED STOCHASTICITY VS. LEARNING THE STOCHASTICITY

CONDITIONS FOR CONVERGENCE

- 1 WEAK GRADIENT DOMINATION $J(\theta^*(\sigma), \sigma) - J(\theta, \sigma) \leq \alpha \|\nabla J(\theta, \sigma)\|_2 + \beta$
- 2 OBJECTIVE SMOOTHNESS $\|\nabla^2 J(\theta, \sigma)\|_2 \leq L_2$
- 3 BOUNDED VARIANCE $\text{Var}[\hat{\nabla} J(\theta, \sigma)] \leq VN^{-1}\sigma^{-2}$

FIXED STOCHASTICITY

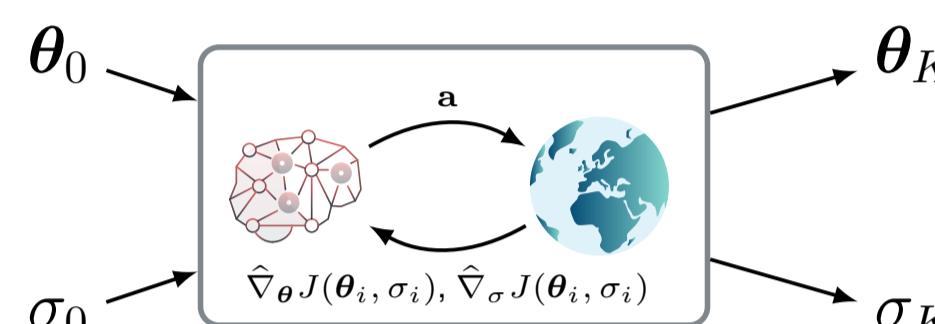


LAST-ITERATE CONVERGENCE

$$\begin{aligned} J(\theta^*(\sigma), \sigma) - \mathbb{E}[J(\theta_K, \sigma)] &\leq \epsilon + \beta \\ \text{WITH } NK &= \tilde{\mathcal{O}}(\epsilon^{-3}\sigma^{-2}) \end{aligned}$$

- ✓ $\sigma = \epsilon \implies$ Deterministic policies in $NK = \tilde{\mathcal{O}}(\epsilon^{-5})$
- ✗ σ has to be selected, it affects $J(\cdot)$

LEARNING STOCHASTICITY



LAST-ITERATE CONVERGENCE

$$\begin{aligned} J^* - \mathbb{E}[J(\theta_K, \sigma_K)] &\leq \epsilon + \beta \\ \text{WITH } NK &= \tilde{\mathcal{O}}(\epsilon^{-3}\sigma_{\min}^{-2}) \end{aligned}$$

- ✓ Select $\bar{\sigma}$, then σ is learned to maximize $J(\cdot)$
- ✗ Final σ is not controllable, no guarantees on deterministic deployment

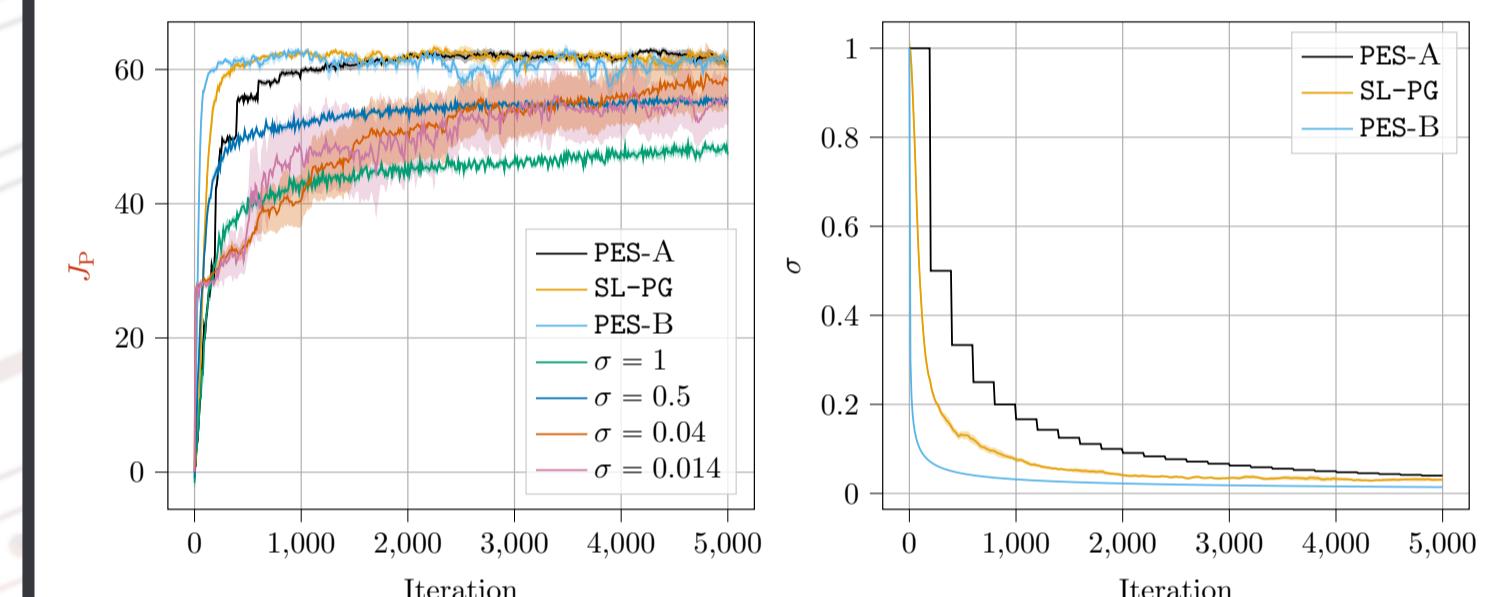
SWITCHING STOCHASTICITY

What happens when the stochasticity σ changes?

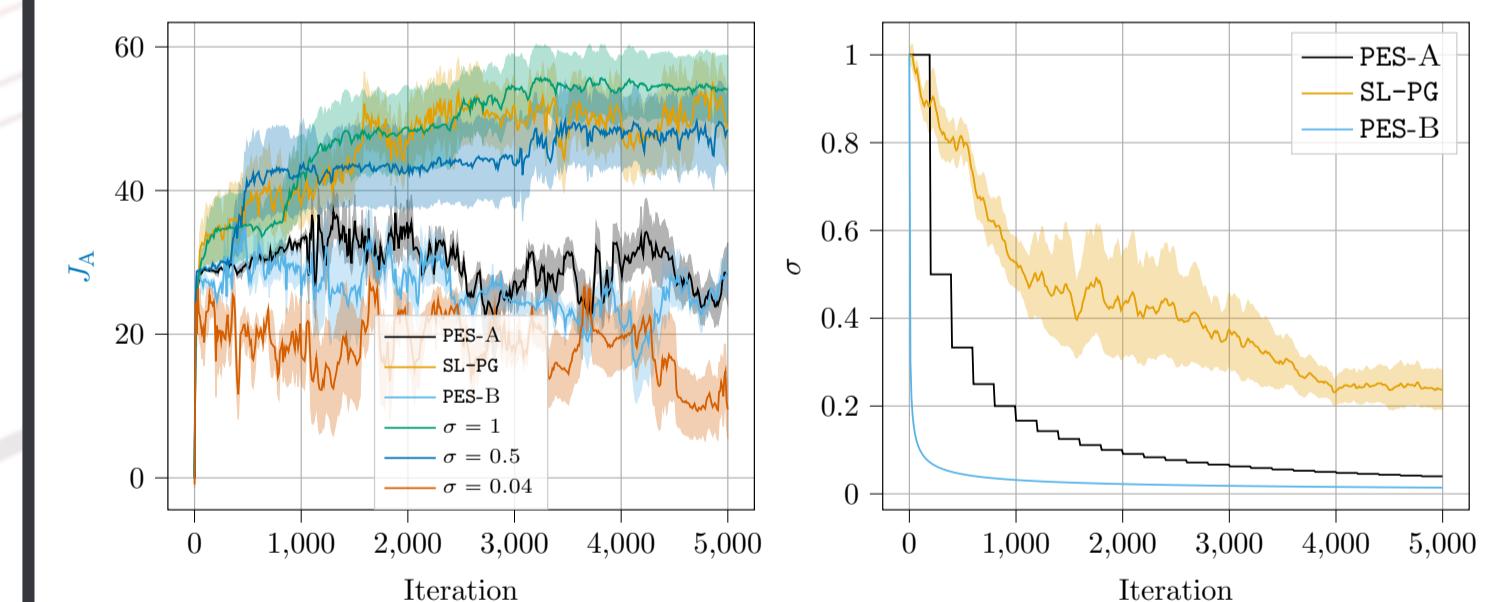
- 1 STOCHASTICITY CHANGE $|J(\theta, \sigma_1) - J(\theta, \sigma_2)| \leq D|\sigma_1 - \sigma_2|$
 - 2 OPTIMUM PLUGGING LOSS $J^*(\sigma_1) - J(\theta^*(\sigma_2), \sigma_1) \leq 2D|\sigma_1 - \sigma_2|$
 - 3 DISTANCE BETWEEN OPTIMA $|J^*(\sigma_1) - J^*(\sigma_2)| \leq D|\sigma_1 - \sigma_2|$
- where $D_A = \mathcal{O}(T\sqrt{d_A})$ and $D_P = \mathcal{O}(\sqrt{d_\Theta})$

DOES IT WORK?

PARAMETER-BASED EXPLORATION



ACTION-BASED EXPLORATION



DETERMINISTIC DEPLOYMENT

Method	PB	Δ_P	AB	Δ_A
PES-A	61.14 ± 1.1	-0.12	28.93 ± 9.89	0.44
SL-PG	61.43 ± 2.73	-0.21	49.83 ± 13.38	-1.19
PES-B	60.40 ± 2.31	-0.27	29.24 ± 8.74	0.33

REFERENCES

- J. Baxter and P. L. Bartlett. Infinite-horizon policy-gradient estimation. *JAIR*, 2001.
A. Montenegro, M. Mussi, A. M. Metelli, and M. Papini. Learning optimal deterministic policies with stochastic policy gradients. *ICML*, 2024.