

LAST-ITERATE GLOBAL CONVERGENCE OF POLICY GRADIENTS FOR CONSTRAINED REINFORCEMENT LEARNING

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



MOTIVATIONS

POLICY GRADIENTS FOR CONSTRAINED RL? PGs handle **continuous state-action spaces**, making them suitable for real-world constrained control problems, learning via **AB** or **PB** exploration paradigms.
CONSTRAINTS ON RISKS? They enforce **safer behaviors** by imposing constraints on risk measures instead of on expected costs.

CONTRIBUTIONS

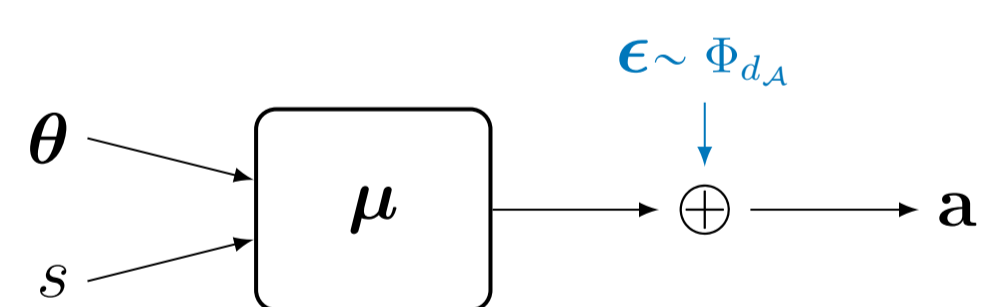
EXPLORATION-AGNOSTIC FRAMEWORK for solving **risk-constrained continuous** control problems via policy-based primal-dual methods under **general policy parameterization** and with either **AB** or **PB** exploration approaches.

LAST-ITERATE GLOBAL CONVERGENCE to a **feasible policy** for the exploration-agnostic method C-PG with:

- **regularized Lagrangian** w.r.t. the dual variable
- **weak gradient domination** w.r.t. the parameters
- **multiple constraints on expected costs**

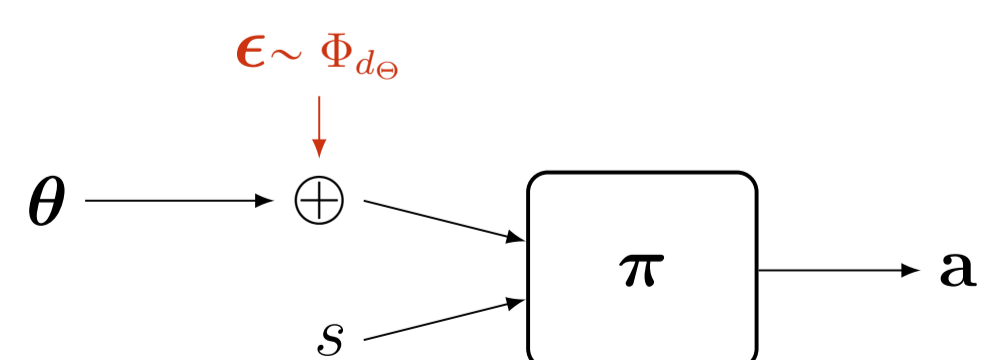
AB AND PB EXPLORATION

ACTION-BASED EXPLORATION



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

PARAMETER-BASED EXPLORATION



$$J_P(\theta) = \mathbb{E}_{\theta \sim \nu_\rho} \left[\mathbb{E}_{\tau \sim p_A(\cdot|\theta)} [R(\tau)] \right]$$

SOLVING RISK CONSTRAINED OPTIMIZATION PROBLEM VIA PGs

EXPLORATION-AGNOSTIC CONSTRAINED OPTIMIZATION PROBLEM

$$\min_{v \in \mathcal{V}} J_0(v) \quad \text{s.t.} \quad J_i(v) \leq b_i, \quad \forall i \in [U]$$

which is equivalent to

$$\min_{v \in \mathcal{V}} \max_{\lambda \geq 0_U} \mathcal{L}_\omega(v, \lambda)$$

ALGORITHM

Projected Alternate Ascent-Descent on the Ridge-Regularized Lagrangian $\mathcal{L}_\omega(v, \lambda)$ w.r.t. λ

$$\mathcal{L}_\omega(v, \lambda) := J_0(v) + \langle \lambda, \mathbf{J}(v) - \mathbf{b} \rangle - \frac{\omega}{2} \|\lambda\|_2^2$$

UNIFIED RISK MEASURE

$$\min_{\eta \in \mathbb{R}} \mathcal{J}(\theta, \eta)$$

→

$$\mathcal{J}(\theta, \eta) := \mathbb{E}_{\tau \sim p(\cdot|\theta)} [f(C(\tau), \eta)] + g(\eta)$$

RISK MAPPING

COST	$C(\tau)$	0
MV	$C(\tau)(1-2\kappa\eta+\kappa C(\tau))$	$\kappa\eta^2$
CVAR $_\alpha$	$(1-\alpha)^{-1}(C(\tau)-\eta)^+$	η
CHANCE	$\mathbb{1}\{C(\tau) \geq n\}$	0

LAST-ITERATE GLOBAL CONVERGENCE

ASSUMPTIONS | ① ψ -GRADIENT DOMINATION | ② \mathcal{L}_ω REGULARITY | ③ $\text{Var}[\widehat{\nabla} \mathcal{L}_\omega]$ BOUNDED | ④ SADDLE POINT EXISTENCE

$$\psi\text{-GRADIENT DOMINATION} \quad \psi \in [1, 2] \quad \rightarrow \quad \|\nabla_v \mathcal{L}_0(v, \lambda)\|_2^\psi \geq \alpha_1 \left(\mathcal{L}_0(v, \lambda) - \min_{v' \in \mathcal{V}} \mathcal{L}_0(v', \lambda) \right) - \beta_1$$

Under 1-4 and with constraints on expected costs

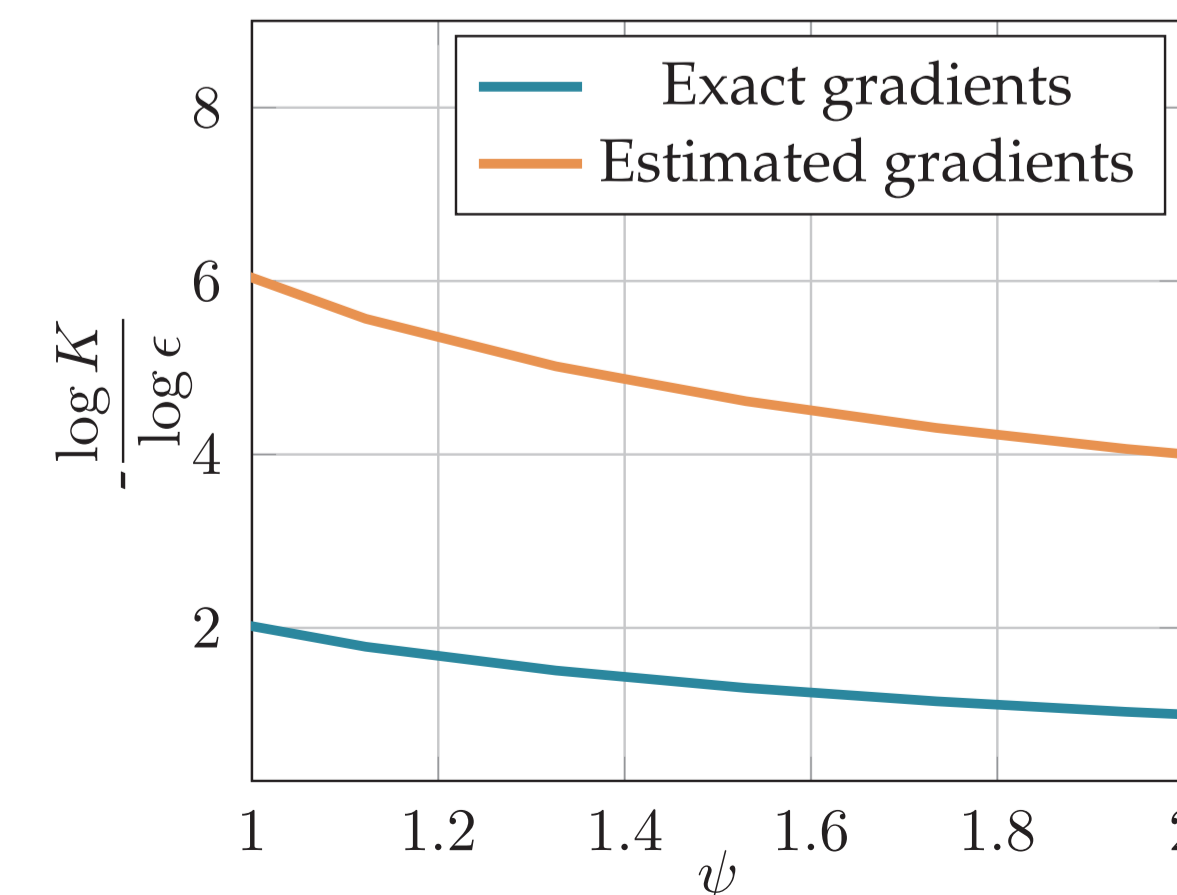
$$\mathbb{E}[J_0(\mathbf{v}_k) - J_0(\mathbf{v}_0^*)] \leq \epsilon + \frac{\beta_1}{\alpha_1} + \frac{\omega}{2} \|\lambda_0^*\|_2^2 \quad \text{and} \quad \mathbb{E}[(J_i(\mathbf{v}_k) - b_i)^+] \leq 4\epsilon + 4\frac{\beta_1}{\alpha_1} + \omega \|\lambda_0^*\|_2, \quad \forall i \in [U]$$

EXACT GRADIENTS

	$\psi = 1$ (GD)	$\psi = 2$ (PL)
FIXED ω	$\omega^{-1}\epsilon^{-1}$	$\omega^{-1}\log(\epsilon^{-1})$
$\omega = \mathcal{O}(\epsilon)$	ϵ^{-2}	$\epsilon^{-1}\log(\epsilon^{-1})$

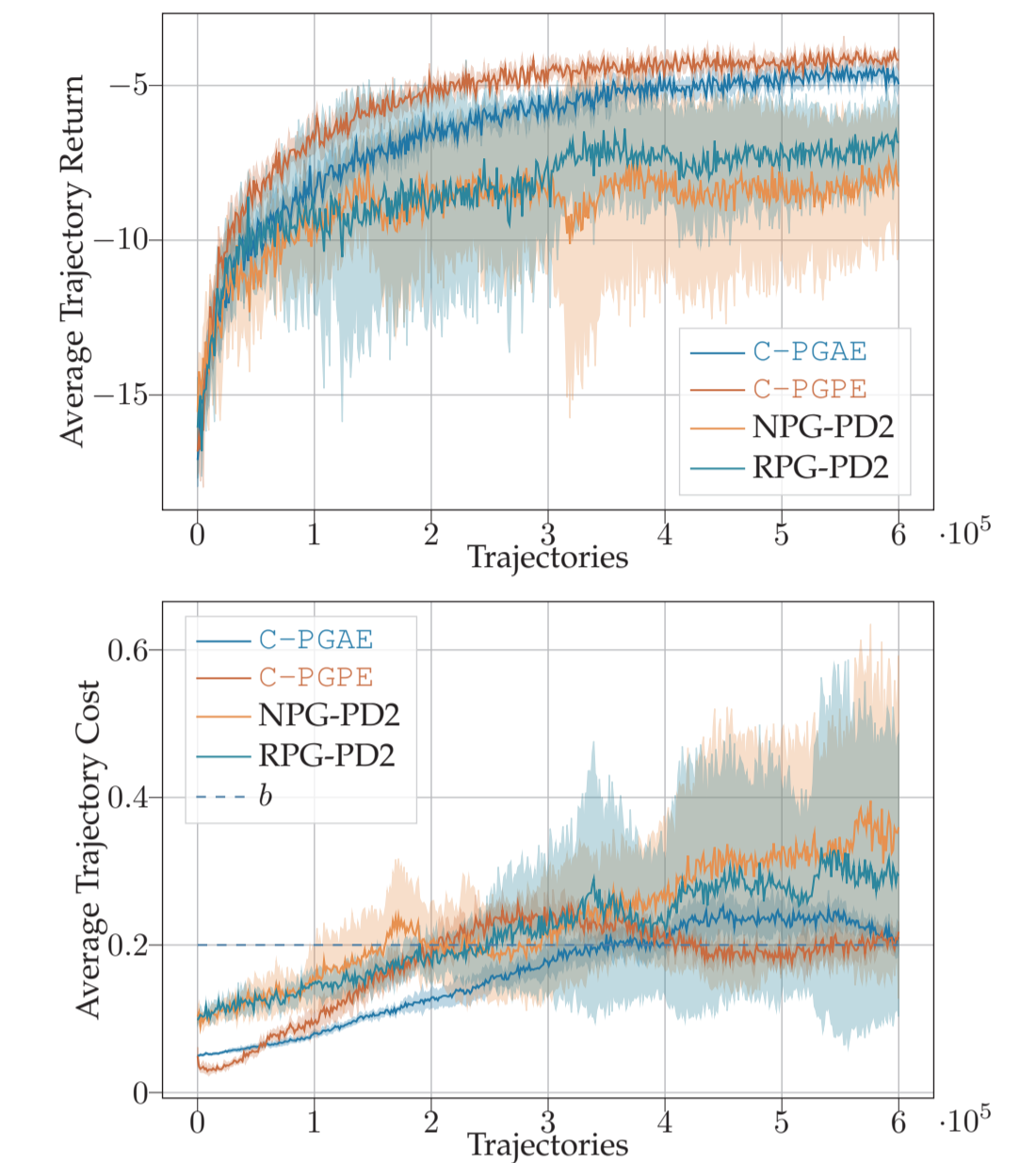
ESTIMATED GRADIENTS

	$\psi = 1$ (GD)	$\psi = 2$ (PL)
FIXED ω	$\omega^{-3}\epsilon^{-3}\log(\epsilon^{-1})$	$\omega^{-3}\epsilon^{-1}\log(\epsilon^{-1})$
$\omega = \mathcal{O}(\epsilon)$	$\epsilon^{-6}\log(\epsilon^{-1})$	$\epsilon^{-4}\log(\epsilon^{-1})$



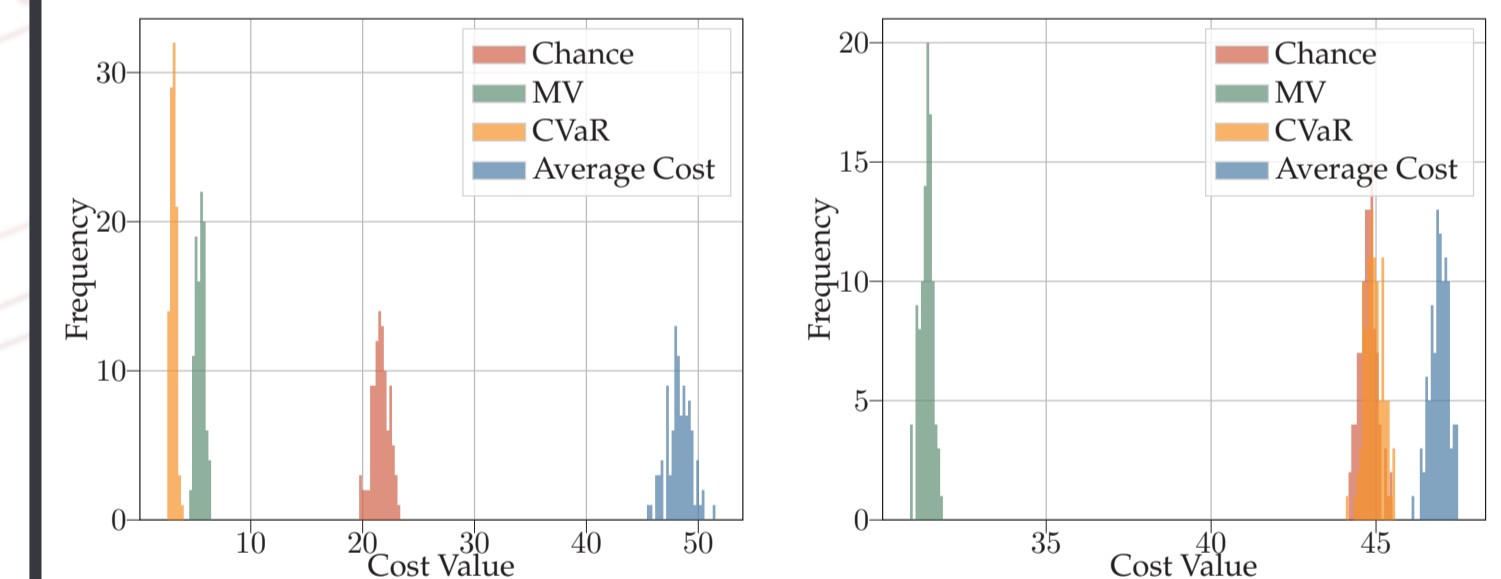
EXPERIMENTS

NUMBER OF TRAJECTORIES STUDY



RISK MINIMIZATION STUDY

LAST-ITERATE COST DISTRIBUTION



LAST-ITERATE PERFORMANCE

COST	CVAR $_\alpha$	MV	CHANCE
26.91 ± 0.09	23.07 ± 0.26	24.23 ± 0.23	26.34 ± 0.16
25.80 ± 0.14	23.08 ± 0.24	0.75 ± 0.26	23.07 ± 0.26

REFERENCES

- Jonathan Baxter and Peter L. Bartlett. Infinite-horizon policy-gradient estimation. *JAIR*, 2001.
- Dongsheng Ding, Chen-Yu Wei, Kaiqing Zhang, and Alejandro Ribeiro. Last-iterate convergent policy gradient primal-dual methods for constrained mdp. *NeurIPS*, 36, 2024.
- Alessandro Montenegro, Marco Mussi, Alberto Maria Metelli, and Matteo Papini. Learning optimal deterministic policies with stochastic policy gradients. *ICML*, 2024.
- Frank Sehnke, Christian Osendorfer, Thomas Rückstieß, Alex Graves, Jan Peters, and Jürgen Schmidhuber. Parameter-exploring policy gradients. *Neural Networks*, 2010.