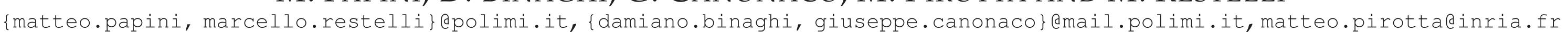


STOCHASTIC VARIANCE-REDUCED POLICY GRADIENT

M. Papini, D. Binaghi, G. Canonaco, M. Pirotta and M. Restelli





MOTIVATION

- ► We want to solve continuous Markov Decision Processes (MDPs), e.g., robot locomotion
- ▶ Policy Gradient (PG): optimize a parametric policy π_{θ} via **gradient ascent** on performance $J(\theta)$.
- ► Two main strategies for gradient computation:
 - Full Gradient (FG): sample infefficient $\longrightarrow O(N/\epsilon)$
 - Stochastic Gradient (SG): slow convergence $\longrightarrow O(1/\epsilon^2)$

Sample efficiency is crucial in Reinforcement Learning (RL), where collecting samples is extremely expensive

⇒ **FG** often unfeasible

Slower convergence of SG due to high gradient variance.

- ► Solution from finite-sum optimization in Supervised Learning (SL):
 - Stochastic Variance-Reduced Gradient (SVRG) $\longrightarrow O(N + N^{2/3}/\epsilon)$

SVRG IN RL: CHALLENGES

Goal: replace SG in RL with SVRG

RL introduces three challenges

- 1. Non-concavity: policy performance $J(\theta)$ is typically a non-concave objective
- 2. Infinite dataset: expectation over all possible trajectories cannot be expressed as a finite sum
- 3. Non-stationarity: the data-generating distribution changes as we learn (policy changes)

WHAT IS OUT THERE?

From the **SL** literature, separate study of:

- Non-concavity [Allen-Zhu and Hazan, 2016, Reddi et al., 2016]
- Infinite dataset [Harikandeh et al., 2015, Bietti and Mairal, 2017]

From the **RL** literature:

- **SVRPGE** [Du et al., 2017]: policy evaluation
- SVRPO [Xu et al., 2017]: direct application of SVRG to TRPO

Non-stationarity never addressed explicitly!

SVRPG: STOCHASTIC VARIANCE-REDUCED POLICY GRADIENT

We design an **SVRG**-like algorithm for the **PG** framework.

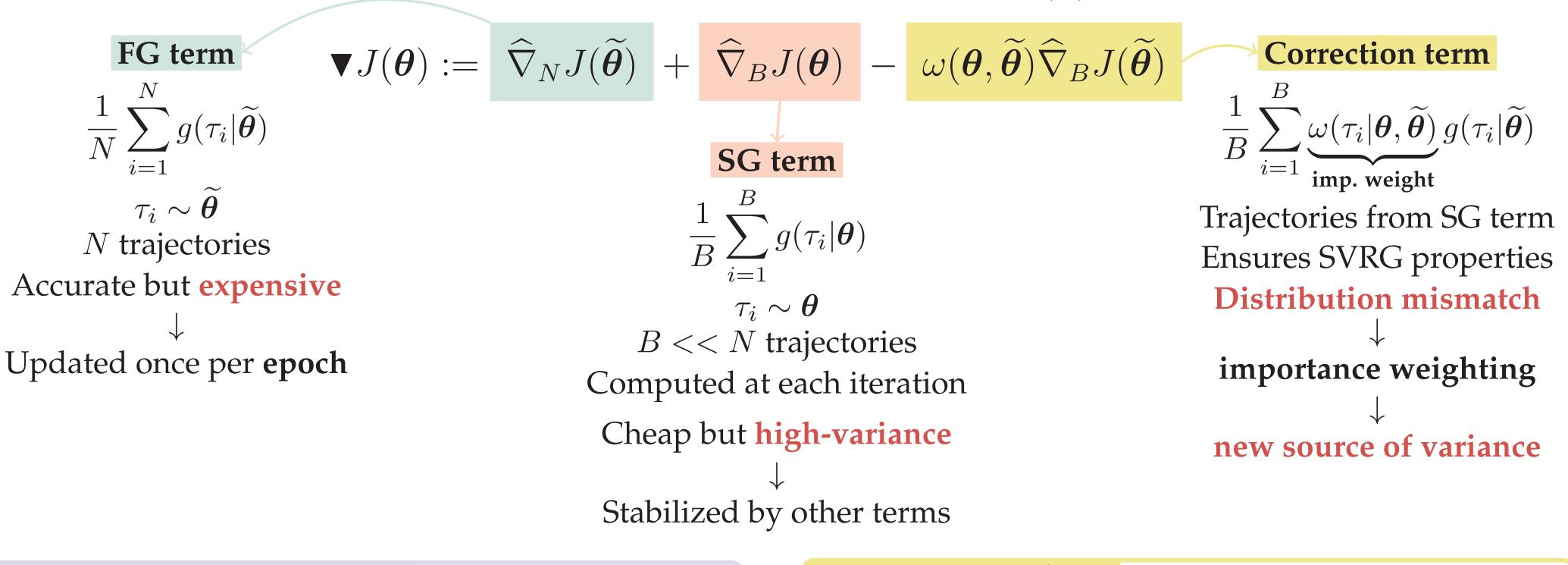
The easy part:

- Take any **unbiased** policy gradient estimator $g(\tau|\boldsymbol{\theta})$ (e.g., REINFORCE)
- SVRG idea: combine frequent SG with rare **FG** into a corrected gradient estimator $\nabla J(\theta)$
- Two-time scales: 1 epoch = m iterations

Solving the RL challenges:

- Manage non-concavity with smooth policies (e.g., Gaussian)
- Approximate infinite dataset with a large batch size
- Correct nonstationarity with importance weighting

Introducing the **SVRPG** estimator $\nabla J(\theta)$:



Fundamental SVRG properties:

- Unbiasedness: $\mathbb{E}\left[lackbox{\psi} J(oldsymbol{ heta})
 ight] =
 abla J(oldsymbol{ heta})$
- 2. Vanishing variance: \mathbb{V} ar $[\nabla J(\theta)] \to 0$ as $\theta \to \theta^*$

Importance weighting

Use samples from θ as taken with θ :

$$w(\tau|\boldsymbol{\theta}, \widetilde{\boldsymbol{\theta}}) = \frac{p(\tau|\boldsymbol{\theta})}{p(\tau|\boldsymbol{\theta})} \quad \leftarrow \text{target distribution} \\ \leftarrow \text{source distribution}$$

ALGORITHM

$$\begin{aligned} & \text{For } s = 1, \dots \\ & \text{Sample } N \text{ trajectories using } \widetilde{\theta} \\ & \text{Compute } \text{FG} = \widehat{\nabla}_N J(\widetilde{\theta}) \\ & \text{For } t = 1, \dots, m \\ & \text{Sample } B \text{ trajectories using } \theta \\ & \text{Compute } \text{SG} = \widehat{\nabla}_B J(\theta) \\ & \text{Compute correction} = \omega(\theta, \widetilde{\theta}) \widehat{\nabla}_B J(\widetilde{\theta}) \\ & \text{Update } \theta \leftarrow \theta + \alpha \blacktriangledown J(\theta) \end{aligned} \end{aligned} \end{aligned}$$

REFERENCES

Zeyuan Allen-Zhu and Elad Hazan. Variance reduction for faster non-convex optimization. In International Conference on Machine Learning, 2016. Alberto Bietti and Julien Mairal. Stochastic optimization with variance reduction for infinite datasets with finite

sum structure. In Advances in Neural Information Processing Systems, 2017. Simon S. Du, Jianshu Chen, Lihong Li, Lin Xiao, and Dengyong Zhou. Stochastic variance reduction methods

for policy evaluation. In ICML, volume 70 of Proceedings of Machine Learning Research, 2017. Reza Harikandeh, Mohamed Osama Ahmed, Alim Virani, Mark Schmidt, Jakub Konečný, and Scott Sallinen. Stopwasting my gradients: Practical svrg. In Advances in Neural Information Processing Systems, 2015.

Sashank J Reddi, Ahmed Hefny, Suvrit Sra, Barnabas Poczos, and Alex Smola. Stochastic variance reduction for nonconvex optimization. In International conference on machine learning, 2016.

Tianbing Xu, Qiang Liu, and Jian Peng. Stochastic variance reduction for policy gradient estimation. CoRR, abs/1710.06034, 2017.

CONVERGENCE

Non-concavity $\longrightarrow local$ maximum Linear rate under increasing batch sizes:

$$\mathbb{E}\left[\|\nabla J(\boldsymbol{\theta})\|^2\right] \le \frac{J(\boldsymbol{\theta}^*) - J(\boldsymbol{\theta}_0)}{\psi T} + \frac{\zeta}{N} + \frac{\xi}{B}$$

Infinite-dataset error Non-stationarity error from FG approximation from importance weighting take large B (still $\ll N$) take *large* N

- Constants (ψ, ζ, ξ) depend only on step size α and epoch size m
- Theory very conservative dictates metaparameters (α, m, N, B) to guarantee convergence

HEURISTICS

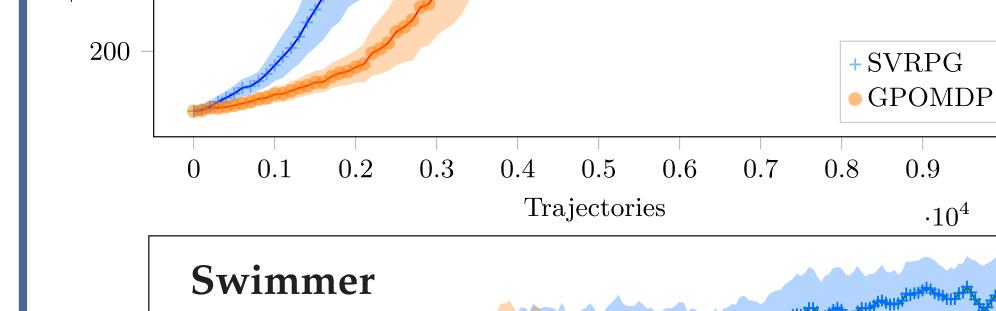
Heuristics can improve performance in practice

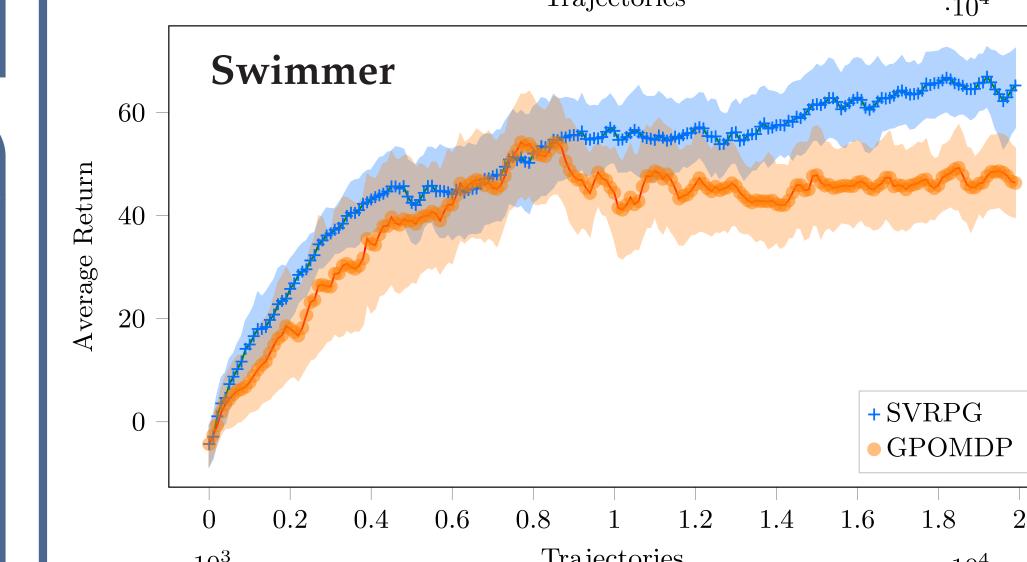
- Step size: ADAM to counteract gradient variance ⚠ FG and SG variance have completely different magnitude → use two separate annealing schedules α_{FG} and α_{SG}
- Epoch size: cut epoch when the effective step size becomes too small

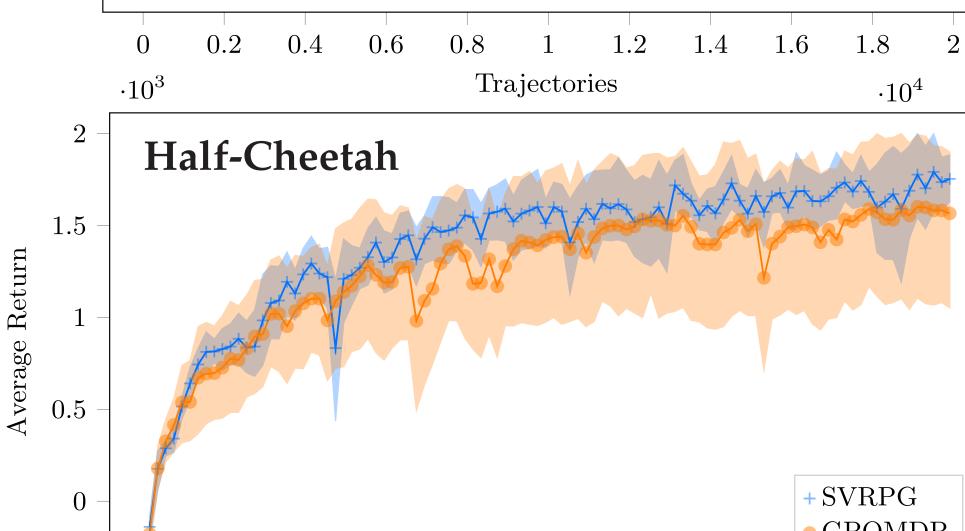
$$\frac{\alpha_{FG}}{N} > \frac{\alpha_{SG}}{B} \implies \text{new epoch}$$

- Normalized importance weights: reduce variance at the price of introducing a small bias
- Critic: an orthogonal variance reduction technique from the **PG** literature

EMPIRICAL RESULTS **Cart-Pole** § 600







Trajectories

