



PROBLEM AND MOTIVATION

- **Reinforcement Learning (RL):** find optimal policy π^*
- **Policy Search:** search over a class of policies π
 - Every policy induces a distribution $p(\cdot|\pi)$ over **trajectories** τ of the Markov Decision Process (MDP)
 - Every trajectory τ has a **return** $R(\tau)$
- **Goal:** find π^* maximizing $J(\pi)$

$$J(\pi) = \mathbb{E}_{\tau \sim p(\cdot|\pi)} [R(\tau)]$$

- Using data collected with some policy π :
 - How can I evaluate proposals $\pi' \neq \pi$?
 - How can I trust counterfactual evaluations?
 - How can I best use my data for optimization?

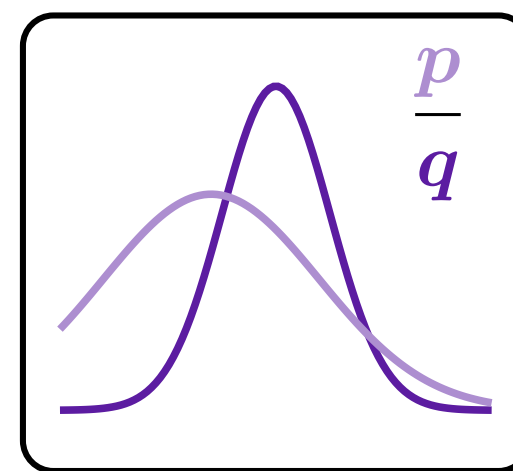
IMPORTANCE SAMPLING

How can I evaluate proposals? With Importance Sampling (IS)

- Given a **behavioral** (data-sampling) distribution $q(x)$, a **target distribution** $p(x)$, and a function $f(x)$, estimate

$$\mu = \mathbb{E}_{x \sim p} [f(x)] \quad \text{with data from } q \quad x_i \sim q$$

$$\hat{\mu}_{\text{IS}} = \frac{1}{N} \sum_{i=1}^N \underbrace{\frac{p(x_i)}{q(x_i)}}_{w(x_i)} f(x_i)$$



- $w(x) = p(x)/q(x)$ is the **importance weight**
- The estimate is **unbiased**: $\mathbb{E}_q[\hat{\mu}_{\text{IS}}] = \mu \dots$
- ... **but the variance can be very high!**
- **Rényi divergence**: dissimilarity between p and q :

$$D_2(p||q) = \log \mathbb{E}_{x \sim q} \left[\left(\frac{p(x)}{q(x)} \right)^2 \right] \quad d_2(p||q) = \exp\{D_2(p||q)\}$$

exponentiated Rényi

- Variance of the weight depends **exponentially** on the distributional divergence (?)

$$\text{Var}[w] = d_2(p||q) - 1$$

- **Effective Sample Size (ESS)**: number of equivalent samples in plain Monte Carlo estimation ($x_i \sim p$)

$$\text{ESS} = \frac{N}{d_2(p||q)} \approx \frac{\|w\|_1^2}{\|w\|_2^2} = \widehat{ESS}$$

- Variance of the estimator $\hat{\mu}_{\text{IS}}$ depends **exponentially** on the distributional divergence as well

$$\text{Var}[\hat{\mu}_{\text{IS}}] \leq \frac{1}{N} \|f\|_\infty^2 d_2(p||q)$$

OFF-DISTRIBUTION LEARNING

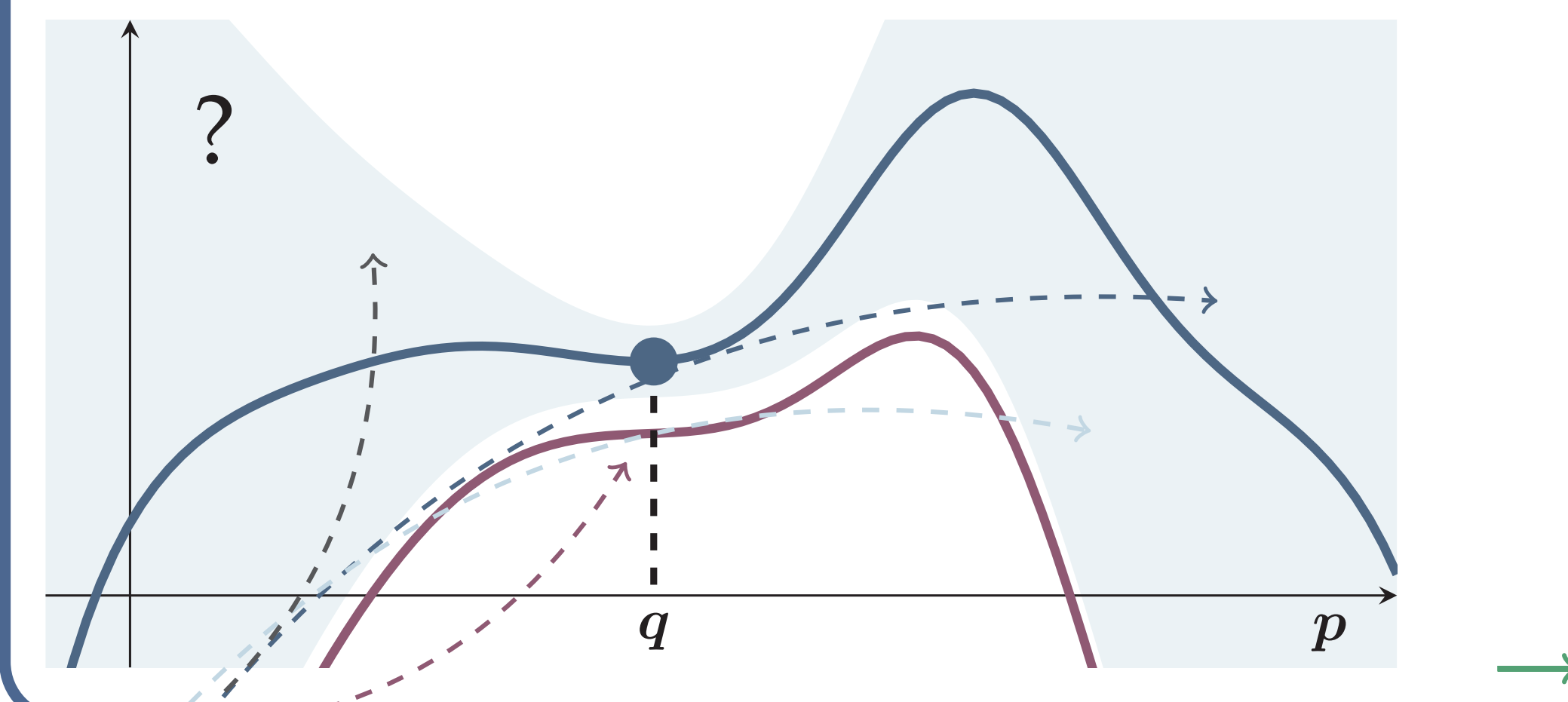
How (far) can I trust counterfactual evaluations?

- Evaluate only close solutions: REPS (?), TRPO (?)
- Use a lower bound: EM (??), PPO (?), POIS

Given a behavioral $q(x)$, a function $f(x)$ and a proposal $p(x)$, with probability at least $1 - \delta$:

$$\mathbb{E}_{x \sim p} [f(x)] \geq \text{Lower Bound} = \frac{1}{N} \sum_{i=1}^N w(x_i) f(x_i) - \|f\|_\infty \sqrt{\frac{(1-\delta)d_2(p||q)}{\delta N}}$$

IS Estimator Variance Bound (Cantelli)



How can I best use my data for optimization?

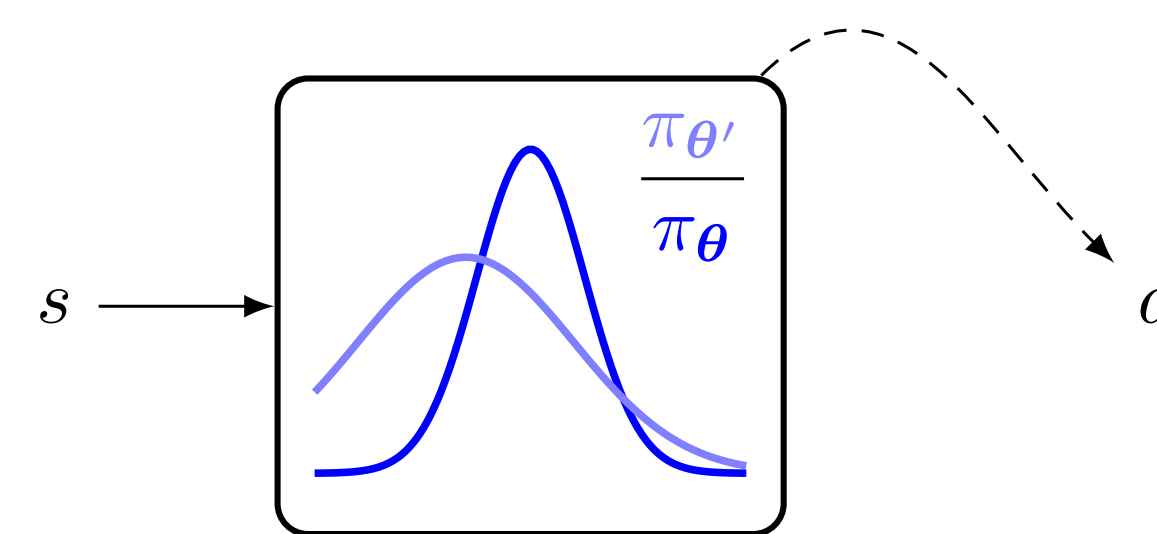
Given the behavioral q , find p maximizing $\mathbb{E}_{x \sim p}[f(x)]$:

1. Collect data with q (expensive in RL)
2. Find p maximizing $\mathcal{L}_\delta^{\text{POIS}}(p/q)$ (offline optimization)
3. Set new behavioral $q \leftarrow p$
4. Repeat until convergence

ACTION-BASED POIS

- Find the **policy** parameters θ^* that maximize $J(\theta')$ (??)

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



- Given a **behavioral policy** π_θ we compute a **target policy** $\pi_{\theta'}$ by optimizing:

$$\mathcal{L}_\lambda^{\text{A-POIS}}(\theta'/\theta) = \frac{1}{N} \sum_{i=1}^N \prod_{t=0}^{H-1} \frac{\pi_{\theta'}(a_{\tau_i,t}|s_{\tau_i,t})}{\pi_\theta(a_{\tau_i,t}|s_{\tau_i,t})} R(\tau_i) - \lambda \sqrt{\frac{\widehat{d}_2(p(\cdot|\theta')||p(\cdot|\theta))}{N}}$$

- The term $d_2(p(\cdot|\theta')||p(\cdot|\theta))$ is estimated from samples
- The d_2 grows exponentially with the task horizon H
- λ is a regularization hyperparameter

$$\lambda = \frac{R_{\max}}{1-\gamma} \sqrt{\frac{1-\delta}{\delta}}$$

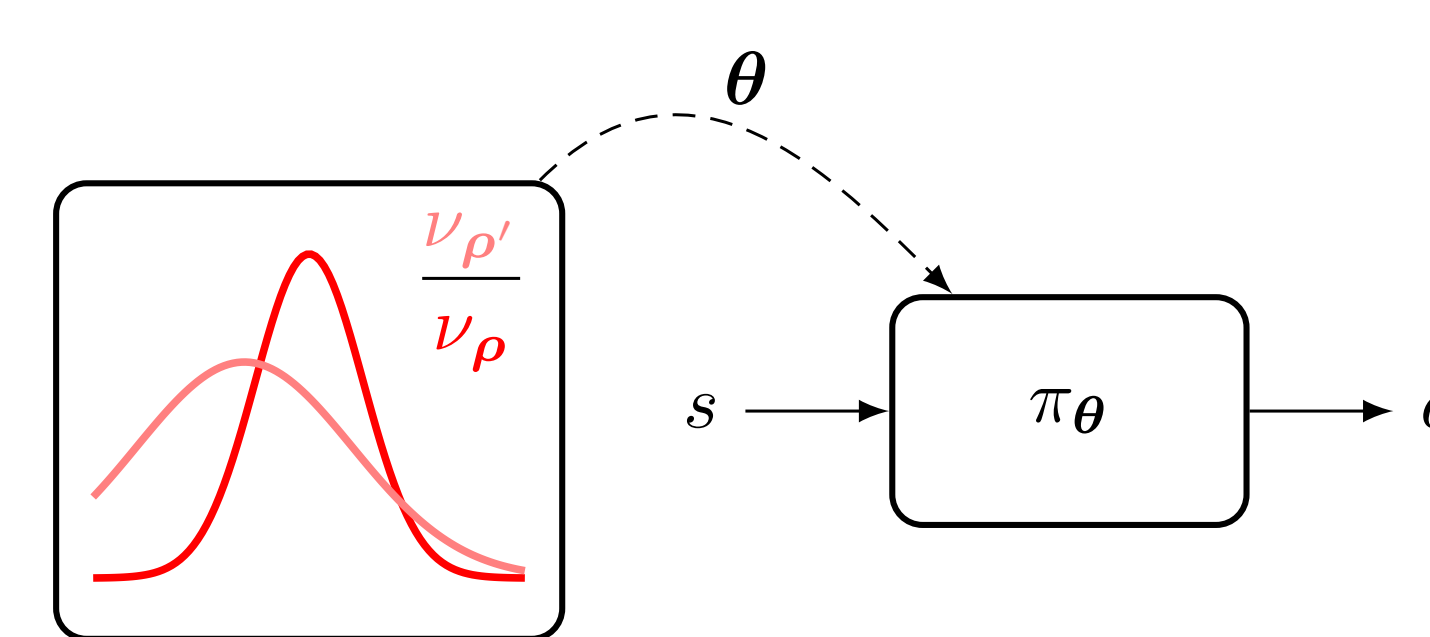
- We consider diagonal Gaussian policies π_θ

$$a \sim \pi_{\mu,\sigma}(\cdot|s) = \mathcal{N}(u_\mu(s), \text{diag}(\sigma^2))$$

PARAMETER-BASED POIS

- Find the **hyperpolicy** parameters ρ^* that maximize $J(\rho)$ (?)

$$J(\rho) = \mathbb{E}_{\theta \sim \nu_\rho} \mathbb{E}_{\tau \sim p(\cdot|\theta)} [R(\tau)]$$



- Given a **behavioral hyperpolicy** ν_ρ we compute a **target hyperpolicy** $\nu_{\rho'}$ by optimizing:

$$\mathcal{L}_\lambda^{\text{P-POIS}}(\rho'/\rho) = \frac{1}{N} \sum_{i=1}^N \frac{\nu_{\rho'}(\theta_i)}{\nu_\rho(\theta_i)} R(\tau_i) - \lambda \sqrt{\frac{d_2(\nu_{\rho'}||\nu_\rho)}{N}}$$

- The term $d_2(\nu_{\rho'}||\nu_\rho)$ can be computed exactly
- Affected by the parameter space dimension $\text{dim}(\theta)$
- λ is a regularization hyperparameter

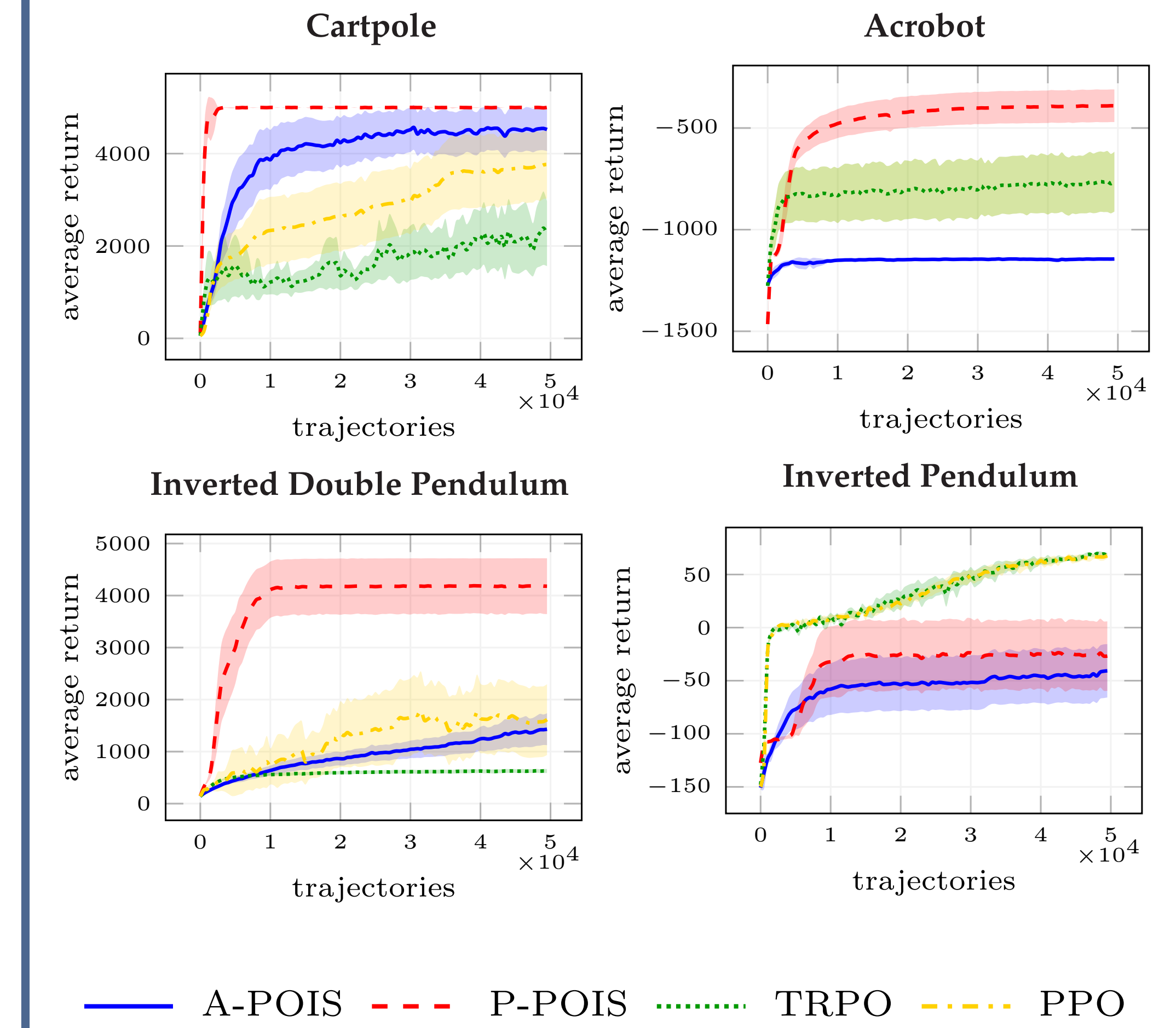
$$\lambda = \frac{R_{\max}}{1-\gamma} \sqrt{\frac{1-\delta}{\delta}}$$

- We consider diagonal Gaussian hyperpolicies ν_ρ

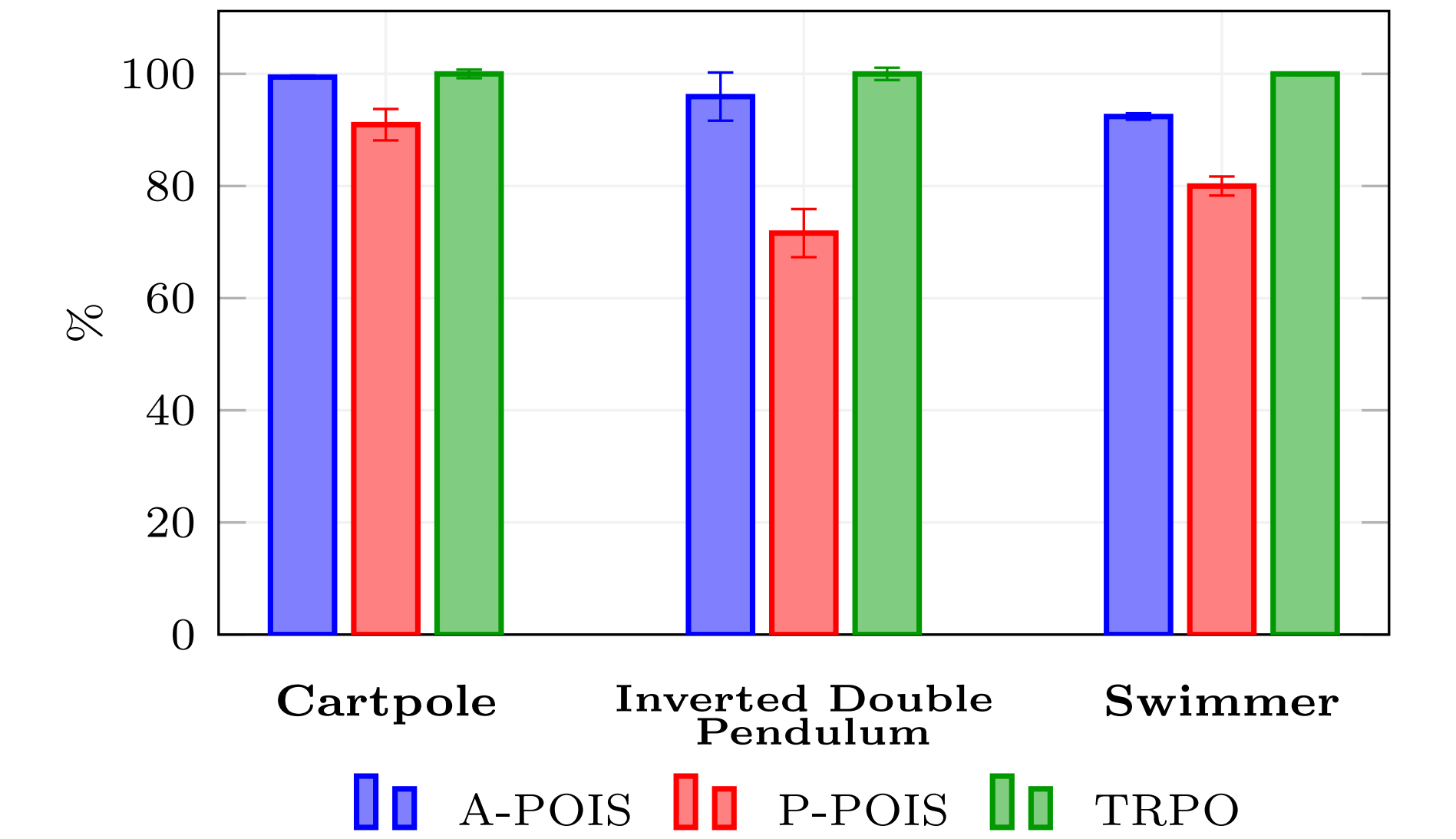
$$\theta \sim \nu_{\mu,\sigma} = \mathcal{N}(\mu, \text{diag}(\sigma^2))$$

EXPERIMENTS

Linear Policies



Deep Policies



Algorithm Details

- **Self-normalized (SN) importance sampling** (?)

$$\tilde{\mu}_{\text{SN}} = \frac{\sum_{i=1}^N w(x_i) f(x_i)}{\sum_{i=1}^N w(x_i)} \quad x_i \sim q$$

- ESS instead of d_2 as penalization
- Gradient optimization of $\mathcal{L}^{\text{A-POIS}}$ using *line search*
- Natural gradient for P-POIS

REFERENCES