

# From Importance Sampling to Doubly Robust Policy Gradient

---



Jiawei Huang (UIUC)



Nan Jiang (UIUC)

Policy Gradient Estimators

Off-Policy Evaluation Estimators

$$\nabla_{\theta} J(\pi_{\theta}) = \lim_{\Delta\theta \rightarrow 0} \frac{J(\pi_{\theta+\Delta\theta}) - J(\pi_{\theta})}{\Delta\theta}$$



Policy Gradient Estimators



Off-Policy Evaluation Estimators

$$\nabla_{\theta} J(\pi_{\theta}) = \lim_{\Delta\theta \rightarrow 0} \frac{J(\pi_{\theta+\Delta\theta}) - J(\pi_{\theta})}{\Delta\theta}$$



Policy Gradient Estimators

REINFORCE

$$\rho_{[0:T]} \sum_{t=0}^T \gamma^t r_t$$



Off-Policy Evaluation Estimators

Traj-wise IS

$$\sum_{t=0}^T \nabla \log \pi_{\theta}^t \sum_{t'=0}^T \gamma^{t'} r_{t'}$$

(Tang and Abbeel, 2010)

Standard PG

$$\sum_{t=0}^T \gamma^t \rho_{[0:t]} r_t$$

Step-wise IS

$$\sum_{t=0}^T \nabla \log \pi_{\theta}^t \sum_{t'=t}^T \gamma^{t'} r_{t'}$$

$$\nabla_{\theta} J(\pi_{\theta}) = \lim_{\Delta\theta \rightarrow 0} \frac{J(\pi_{\theta+\Delta\theta}) - J(\pi_{\theta})}{\Delta\theta}$$



Policy Gradient Estimators



Off-Policy Evaluation Estimators

PG with State Baselines

$$\sum_{t=0}^T \nabla \log \pi_{\theta}^t \left( \sum_{t'=t}^T \gamma^{t'} r_{t'} - \gamma^t b_t \right)$$

OPE with State Baselines

$$b_0 + \sum_{t=0}^T \gamma^t \rho_{[0:t]} \left( r_t + \gamma b_{t+1} - b_t \right)$$

$$\nabla_{\theta} J(\pi_{\theta}) = \lim_{\Delta\theta \rightarrow 0} \frac{J(\pi_{\theta+\Delta\theta}) - J(\pi_{\theta})}{\Delta\theta}$$



Policy Gradient Estimators



Off-Policy Evaluation Estimators

Trajectory-wise CV (Cheng et al., 2019)

$$\sum_{t=0}^T \left\{ \nabla \log \pi_{\theta}^t \left[ \sum_{t'=t}^T \gamma^{t'} r_{t'} + \sum_{t'=t+1}^T \gamma^{t'} (\tilde{V}_{t'}^{\pi_{\theta}} - \tilde{Q}_{t'}^{\pi_{\theta}}) \right] + \gamma^t (\nabla \tilde{V}_t^{\pi_{\theta}} - \tilde{Q}_t^{\pi_{\theta}} \nabla \log \pi_{\theta}^t) \right\}$$

**DR-PG (Ours)**

$$\sum_{t=0}^T \left\{ \nabla \log \pi_{\theta}^t \left[ \sum_{t'=t}^T \gamma^{t'} r_{t'} + \sum_{t'=t+1}^T \gamma^{t'} (\tilde{V}_{t'}^{\pi_{\theta}} - \tilde{Q}_{t'}^{\pi_{\theta}}) \right] + \gamma^t (\nabla \tilde{V}_t^{\pi_{\theta}} - \nabla_{\theta} \tilde{Q}_t^{\pi_{\theta}} - \tilde{Q}_t^{\pi_{\theta}} \nabla \log \pi_{\theta}^t) \right\}$$



Doubly Robust OPE

$$\tilde{V}_0^{\pi'} + \sum_{t=0}^T \gamma^t \rho_{[0:t]} (r_t + \gamma \tilde{V}_{t+1}^{\pi'} - \tilde{Q}_t^{\pi'})$$



## MDP Setting

- Episodic RL with discount factor  $\gamma$ , and maximum episode length  $T$ ;
- Fixed initial state distribution;
- Trajectory is defined as  $s_0, a_0, r_0, s_1, \dots, s_T, a_T, r_T$ .

## Frequently used notations

- $\pi_\theta$ : Policy parameterized by  $\theta$ .
- $J(\pi_\theta) = \mathbb{E}_{\pi_\theta} [\sum_{t=0}^T \gamma^t r(s_t, a_t)]$ : Expected discounted return of  $\pi_\theta$ .

## From Stepwise IS OPE to Standard PG

$\pi_\theta$  is the behavior policy and  $\pi_{\theta+\Delta\theta}$  as the target policy.  $r_t = r(s_t, a_t)$  and  $\pi_\theta^t = \pi_\theta(a_t|s_t)$ .

$$\widehat{J}(\pi_{\theta+\Delta\theta}) = \sum_{t=0}^T \gamma^t r_t \prod_{t'=0}^t \frac{\pi_{\theta+\Delta\theta}^{t'}}{\pi_\theta^{t'}}$$



## From Stepwise IS OPE to Standard PG

$\pi_\theta$  is the behavior policy and  $\pi_{\theta+\Delta\theta}$  as the target policy.  $r_t = r(s_t, a_t)$  and  $\pi_\theta^t = \pi_\theta(a_t|s_t)$ .

$$\begin{aligned}\widehat{J}(\pi_{\theta+\Delta\theta}) &= \sum_{t=0}^T \gamma^t r_t \prod_{t'=0}^t \frac{\pi_{\theta+\Delta\theta}^{t'}}{\pi_\theta^{t'}} \\ &= \sum_{t=0}^T \gamma^t r_t \left( 1 + \sum_{t'=0}^t \frac{\nabla_{\theta} \pi_{\theta}^{t'}}{\pi_{\theta}^{t'}} \right) \Delta\theta + o(\Delta\theta)\end{aligned}$$

## From Stepwise IS OPE to Standard PG

$\pi_\theta$  is the behavior policy and  $\pi_{\theta+\Delta\theta}$  as the target policy.  $r_t = r(s_t, a_t)$  and  $\pi_\theta^t = \pi_\theta(a_t|s_t)$ .

$$\begin{aligned}\widehat{J}(\pi_{\theta+\Delta\theta}) &= \sum_{t=0}^T \gamma^t r_t \prod_{t'=0}^t \frac{\pi_{\theta+\Delta\theta}^{t'}}{\pi_\theta^{t'}} \\ &= \sum_{t=0}^T \gamma^t r_t \left( 1 + \sum_{t'=0}^t \frac{\nabla_{\theta} \pi_{\theta}^{t'}}{\pi_{\theta}^{t'}} \right) \Delta\theta + o(\Delta\theta) \\ &= \widehat{J}(\pi_\theta) + \left( \sum_{t=0}^T \gamma^t r_t \sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}^{t'} \right) \Delta\theta + o(\Delta\theta).\end{aligned}$$

## From Stepwise IS OPE to Standard PG

$\pi_\theta$  is the behavior policy and  $\pi_{\theta+\Delta\theta}$  as the target policy.  $r_t = r(s_t, a_t)$  and  $\pi_\theta^t = \pi_\theta(a_t|s_t)$ .

$$\begin{aligned}\widehat{J}(\pi_{\theta+\Delta\theta}) &= \sum_{t=0}^T \gamma^t r_t \prod_{t'=0}^t \frac{\pi_{\theta+\Delta\theta}^{t'}}{\pi_\theta^{t'}} \\ &= \sum_{t=0}^T \gamma^t r_t \left( 1 + \sum_{t'=0}^t \frac{\nabla_\theta \pi_\theta^{t'}}{\pi_\theta^{t'}} \right) \Delta\theta + o(\Delta\theta) \\ &= \widehat{J}(\pi_\theta) + \left( \sum_{t=0}^T \gamma^t r_t \sum_{t'=0}^t \nabla_\theta \log \pi_\theta^{t'} \right) \Delta\theta + o(\Delta\theta).\end{aligned}$$

Then

$$\lim_{\Delta\theta \rightarrow 0} \frac{\widehat{J}(\pi_{\theta+\Delta\theta}) - \widehat{J}(\pi_\theta)}{\Delta\theta} = \sum_{t=0}^T \gamma^t r_t \sum_{t'=0}^t \nabla_\theta \log \pi_\theta^{t'}$$

which is known to be the standard PG.

**Definition:** Doubly-robust OPE estimator (**unbiased**) (Jiang and Li, 2016)

$$\widehat{J}(\pi_{\theta+\Delta\theta}) = \widetilde{V}_0^{\pi_{\theta+\Delta\theta}} + \sum_{t=0}^T \gamma^t \left( \prod_{t'=0}^t \frac{\pi_{\theta+\Delta\theta}^{t'}}{\pi_{\theta}^{t'}} \right) \left( r_t + \gamma \widetilde{V}_{t+1}^{\pi_{\theta+\Delta\theta}} - \widetilde{Q}_t^{\pi_{\theta+\Delta\theta}} \right).$$

where  $\widetilde{V}^{\theta+\Delta\theta} = \mathbb{E}_{a \sim \pi_{\theta+\Delta\theta}} [\widetilde{Q}^{\theta+\Delta\theta}]$ .

# Doubly-Robust Policy Gradient (DR-PG)

**Definition:** Doubly-robust OPE estimator (**unbiased**) (Jiang and Li, 2016)

$$\widehat{J}(\pi_{\theta+\Delta\theta}) = \widetilde{V}_0^{\pi_{\theta+\Delta\theta}} + \sum_{t=0}^T \gamma^t \left( \prod_{t'=0}^t \frac{\pi_{\theta+\Delta\theta}^{t'}}{\pi_{\theta}^{t'}} \right) \left( r_t + \gamma \widetilde{V}_{t+1}^{\pi_{\theta+\Delta\theta}} - \widetilde{Q}_t^{\pi_{\theta+\Delta\theta}} \right).$$

where  $\widetilde{V}^{\theta+\Delta\theta} = \mathbb{E}_{a \sim \pi_{\theta+\Delta\theta}} [\widetilde{Q}^{\theta+\Delta\theta}]$ .

**Theorem:** Given DR-OPE estimator above, we can derive two **unbiased** estimators:

- If  $\widetilde{Q}^{\pi_{\theta+\Delta\theta}} = \widetilde{Q}^{\pi_{\theta}}$  for arbitrary  $\Delta\theta$  [**Traj-CV**, (Cheng, Yan, and Boots., 2019)]

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} + \sum_{t_2=t+1}^T \gamma^{t_2} \left( \widetilde{V}_{t_2}^{\pi_{\theta}} - \widetilde{Q}_{t_2}^{\pi_{\theta}} \right) \right] + \gamma^t \left( \nabla_{\theta} \widetilde{V}_t^{\pi_{\theta}} - \widetilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t \right) \right\}.$$

- else [**DR-PG (Ours)**]

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} + \sum_{t_2=t+1}^T \gamma^{t_2} \left( \widetilde{V}_{t_2}^{\pi_{\theta}} - \widetilde{Q}_{t_2}^{\pi_{\theta}} \right) \right] + \gamma^t \left( \nabla_{\theta} \widetilde{V}_t^{\pi_{\theta}} - \nabla_{\theta} \widetilde{Q}_t^{\pi_{\theta}} - \widetilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t \right) \right\}.$$

# Doubly-Robust Policy Gradient (DR-PG)

**Definition:** Doubly-robust OPE estimator (Jiang and Li, 2016)

$$\widehat{J}(\pi_{\theta+\Delta\theta}) = \widetilde{V}_0^{\pi_{\theta+\Delta\theta}} + \sum_{t=0}^T \gamma^t \left( \prod_{t'=0}^t \frac{\pi_{\theta+\Delta\theta}^{t'}}{\pi_{\theta}^{t'}} \right) \left( r_t + \gamma \widetilde{V}_{t+1}^{\pi_{\theta+\Delta\theta}} - \widetilde{Q}_t^{\pi_{\theta+\Delta\theta}} \right).$$

**Theorem:** Given DR-OPE estimator above, we can derive:

- If  $\widetilde{Q}^{\pi_{\theta+\Delta\theta}} = \widetilde{Q}^{\pi_{\theta}}$  for arbitrary  $\Delta\theta$  [Traj-CV, (Cheng, Yan, and Boots., 2019)]

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} + \sum_{t_2=t+1}^T \gamma^{t_2} \left( \widetilde{V}_{t_2}^{\pi_{\theta}} - \widetilde{Q}_{t_2}^{\pi_{\theta}} \right) \right] + \gamma^t \left( \nabla_{\theta} \widetilde{V}_t^{\pi_{\theta}} - \widetilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t \right) \right\}.$$

- else [DR-PG (Ours)]

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} + \sum_{t_2=t+1}^T \gamma^{t_2} \left( \widetilde{V}_{t_2}^{\pi_{\theta}} - \widetilde{Q}_{t_2}^{\pi_{\theta}} \right) \right] + \gamma^t \left( \nabla_{\theta} \widetilde{V}_t^{\pi_{\theta}} - \nabla_{\theta} \widetilde{Q}_t^{\pi_{\theta}} - \widetilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t \right) \right\}.$$

**Remark 1:** The definitions of  $\nabla_{\theta} \widetilde{V}$  are different. In Traj-CV,  $\nabla_{\theta} \widetilde{V} = \mathbb{E}_{\pi_{\theta}} [\widetilde{Q}^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}]$ , while in DR-PG,  $\nabla_{\theta} \widetilde{V} = \mathbb{E}_{\pi_{\theta}} [\widetilde{Q}^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta} + \nabla_{\theta} \widetilde{Q}^{\pi_{\theta}}]$

**Remark 2:**  $\nabla_{\theta} \widetilde{Q}^{\pi_{\theta}}$  is not necessary a gradient but just an approximation of  $\nabla_{\theta} Q^{\pi_{\theta}}$ .

## DR-PG

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} + \sum_{t_2=t+1}^T \gamma^{t_2} (\tilde{V}_{t_2}^{\pi_{\theta}} - \tilde{Q}_{t_2}^{\pi_{\theta}}) \right] + \gamma^t (\nabla_{\theta} \tilde{V}_t^{\pi_{\theta}} - \nabla_{\theta} \tilde{Q}_t^{\pi_{\theta}} - \tilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t) \right\}.$$

## DR-PG

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} + \sum_{t_2=t+1}^T \gamma^{t_2} (\tilde{V}_{t_2}^{\pi_{\theta}} - \tilde{Q}_{t_2}^{\pi_{\theta}}) \right] + \gamma^t (\nabla_{\theta} \tilde{V}_t^{\pi_{\theta}} - \nabla_{\theta} \tilde{Q}_t^{\pi_{\theta}} - \tilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t) \right\}.$$

Use  $\tilde{Q}^{\pi'}$  invariant to  $\pi'$   $\downarrow$  Traj-CV

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} + \sum_{t_2=t+1}^T \gamma^{t_2} (\tilde{V}_{t_2}^{\pi_{\theta}} - \tilde{Q}_{t_2}^{\pi_{\theta}}) \right] + \gamma^t (\nabla_{\theta} \tilde{V}_t^{\pi_{\theta}} - \tilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t) \right\}.$$



## DR-PG

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} + \sum_{t_2=t+1}^T \gamma^{t_2} (\tilde{V}_{t_2}^{\pi_{\theta}} - \tilde{Q}_{t_2}^{\pi_{\theta}}) \right] + \gamma^t (\nabla_{\theta} \tilde{V}_t^{\pi_{\theta}} - \nabla_{\theta} \tilde{Q}_t^{\pi_{\theta}} - \tilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t) \right\}.$$

Use  $\tilde{Q}^{\pi'}$  invariant to  $\pi'$   $\downarrow$  Traj-CV

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} + \sum_{t_2=t+1}^T \gamma^{t_2} (\tilde{V}_{t_2}^{\pi_{\theta}} - \tilde{Q}_{t_2}^{\pi_{\theta}}) \right] + \gamma^t (\nabla_{\theta} \tilde{V}_t^{\pi_{\theta}} - \tilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t) \right\}.$$

$\mathbb{E} \left[ \sum_{t_2=t+1}^T \gamma^{t_2} (\tilde{V}_{t_2}^{\pi_{\theta}} - \tilde{Q}_{t_2}^{\pi_{\theta}}) \middle| S_{t+1} \right] = 0$ , dropped  $\downarrow$  PG with state-action baselines

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} \right] + \gamma^t (\nabla_{\theta} \tilde{V}_t^{\pi_{\theta}} - \tilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t) \right\}.$$

## DR-PG

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} + \sum_{t_2=t+1}^T \gamma^{t_2} (\tilde{V}_{t_2}^{\pi_{\theta}} - \tilde{Q}_{t_2}^{\pi_{\theta}}) \right] + \gamma^t (\nabla_{\theta} \tilde{V}_t^{\pi_{\theta}} - \nabla_{\theta} \tilde{Q}_t^{\pi_{\theta}} - \tilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t) \right\}.$$

Use  $\tilde{Q}^{\pi'}$  invariant to  $\pi'$   $\downarrow$  Traj-CV

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} + \sum_{t_2=t+1}^T \gamma^{t_2} (\tilde{V}_{t_2}^{\pi_{\theta}} - \tilde{Q}_{t_2}^{\pi_{\theta}}) \right] + \gamma^t (\nabla_{\theta} \tilde{V}_t^{\pi_{\theta}} - \tilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t) \right\}.$$

$\mathbb{E} \left[ \sum_{t_2=t+1}^T \gamma^{t_2} (\tilde{V}_{t_2}^{\pi_{\theta}} - \tilde{Q}_{t_2}^{\pi_{\theta}}) \middle| s_{t+1} \right] = 0$ , dropped  $\downarrow$  PG with state-action baselines

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} \right] + \gamma^t (\nabla_{\theta} \tilde{V}_t^{\pi_{\theta}} - \tilde{Q}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t) \right\}.$$

Use  $\tilde{V}$  as  $\tilde{Q}$   $\downarrow$  PG with state baselines

$$\sum_{t=0}^T \left\{ \nabla_{\theta} \log \pi_{\theta}^t \left[ \sum_{t_1=t}^T \gamma^{t_1} r_{t_1} \right] + \gamma^t (-\tilde{V}_t^{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}^t) \right\}.$$

**Theorem** The covariance matrix of the DR-PG estimator is

$$\begin{aligned}
 & \mathbb{E} \left[ \underbrace{\sum_{n=0}^T \gamma^{2n} \left( \mathbb{V}_{n+1}[r_n] \left( \sum_{t=0}^n \nabla_{\theta} \log \pi_{\theta}^t \right) \left( \sum_{t=0}^n \nabla_{\theta} \log \pi_{\theta}^t \right)^{\top} \right)}_{\text{Randomness of reward}} \right. \\
 & \quad + \underbrace{\text{Cov}_n \left[ \nabla_{\theta} V_n^{\pi_{\theta}} + \left( \sum_{t=0}^{n-1} \nabla_{\theta} \log \pi_{\theta}^t \right) V_n^{\pi_{\theta}} \right]}_{\text{Randomness of transition}} \\
 & \quad \left. + \underbrace{\text{Cov}_n \left[ \nabla_{\theta} Q_n^{\pi_{\theta}} - \nabla_{\theta} \tilde{Q}_n^{\pi_{\theta}} + \left( \sum_{t=0}^n \nabla_{\theta} \log \pi_{\theta}^t \right) \left( Q_n^{\pi_{\theta}} - \tilde{Q}_n^{\pi_{\theta}} \right) \middle| S_n \right]}_{\text{Randomness of policy}} \right].
 \end{aligned}$$

where

$$\begin{aligned}
 \mathbb{V}_n[\cdot] & := \mathbb{V}[\cdot | S_0, a_0, \dots, S_{n-1}, a_{n-1}] \\
 \mathbb{E}_n[\cdot] & := \mathbb{E}[\cdot | S_0, a_0, \dots, S_{n-1}, a_{n-1}] \\
 \text{Cov}_n[\mathbf{v}] & := \mathbb{E}_n[\mathbf{v}\mathbf{v}^{\top}] - \mathbb{E}_n[\mathbf{v}]\mathbb{E}_n[\mathbf{v}]^{\top}.
 \end{aligned}$$

**Theorem:** For tree-structured MDPs (i.e., each state only appears at a unique time step and can be reached by a unique trajectory), the Cramer-Rao lower bound of PG is

$$\mathbb{E} \left[ \sum_{t=0}^T \gamma^{2t} \left\{ \underbrace{\mathbb{V}_{t+1}[r_t] \left[ \left( \sum_{t_1=0}^t \frac{\partial \log \pi_{\theta}^{t_1}}{\partial \theta_i} \right) \right]^2}_{\text{Randomness of reward}} + \underbrace{\mathbb{V}_t \left[ \left( V_t^{\pi_{\theta}} \sum_{t_1=0}^{t-1} \frac{\partial \log \pi_{\theta}^{t_1}}{\partial \theta_i} + \frac{\partial V_t^{\pi_{\theta}}}{\partial \theta_i} \right) \right]}_{\text{Randomness of Transition}} \right\} \right],$$

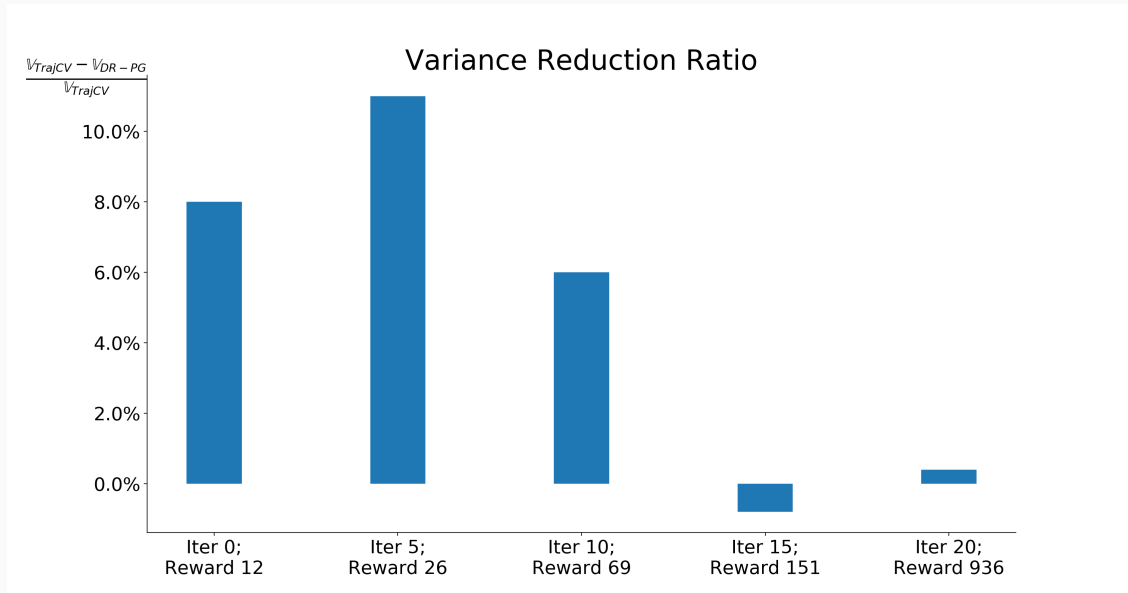
which coincides with the variance of DR-PG when  $\tilde{Q}^{\pi_{\theta}} \equiv Q^{\pi_{\theta}}$  and  $\nabla_{\theta} \tilde{Q}^{\pi_{\theta}} \equiv \nabla_{\theta} Q^{\pi_{\theta}}$ .

## Covariance Comparison in Special Case

Deterministic environment with perfect value function estimation

Estimator	Covariance Matrices
PG with state baselines	$\mathbb{E} \left[ \sum_n \text{Cov}_n \left[ \nabla_\theta Q_n^{\pi_\theta} + \left( \sum_{t=0}^{n-1} \nabla_\theta \log \pi_\theta^t \right) Q_n^{\pi_\theta} + \nabla_\theta \log \pi_\theta^n A_n^{\pi_\theta} \mid S_n \right] \right]$
PG with state-action baselines	$\mathbb{E} \left[ \sum_n \text{Cov}_n \left[ \nabla_\theta Q_n^{\pi_\theta} + \left( \sum_{t=0}^{n-1} \nabla_\theta \log \pi_\theta^t \right) Q_n^{\pi_\theta} \mid S_n \right] \right]$
Trajwise-CV	$\mathbb{E} \left[ \sum_n \text{Cov}_n \left[ \nabla_\theta Q_n^{\pi_\theta} \mid S_n \right] \right]$
DR-PG	0

# Experiments (Variance Reduction)



**Figure 1:** Variance reduction ratio.  $\nabla_G$  denotes the sum of estimator  $G$ 's variance over all parameters of the neural network.

# Experiments (Algorithm Performance)

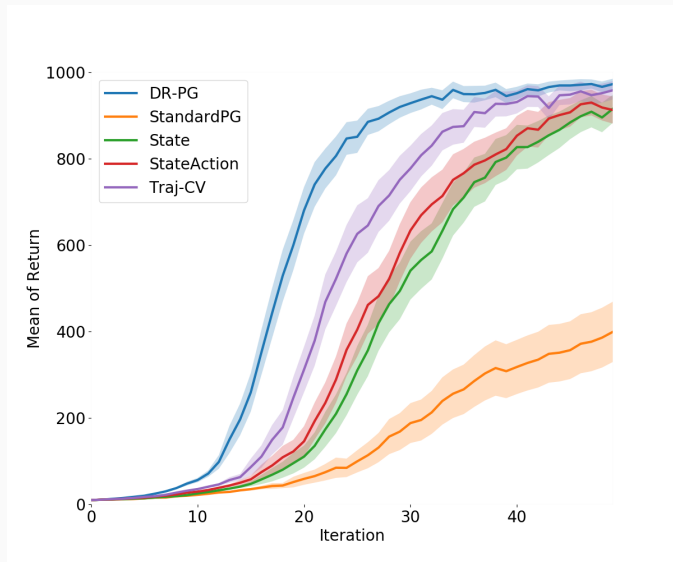


Figure 2: Performance in CartPole task. Average over 150 trials. Plot twice standard error.

# Experiments (Algorithm Time Complexity)

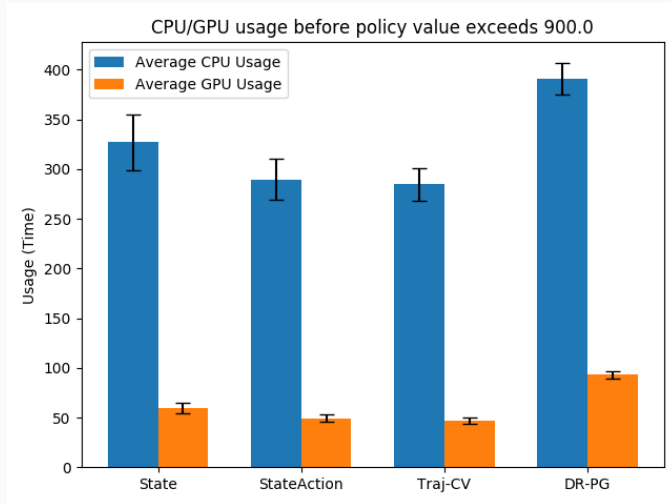


Figure 3: Comparison of GPU/CPU Usage .



Thank You!

Welcome to our Q&A session!

---