
Supplementary Material to Trajectory-Based Off-Policy Deep RL

Anonymous Authors¹

A. Proof of Proposition 1

The weighted importance sampling estimator of the expected cost is given by

$$\hat{j}^{\text{WIS}}(\theta) = \frac{1}{\sum_{i=0}^M w(\tau_i, \theta)} \sum_{i=0}^M w(\tau_i, \theta) R(\tau_i), \quad (1)$$

as derived in Sec. 3. Taking the derivative with respect to the policy parameters, we obtain the policy gradient formulation from theorem 1 as shown in (7).

B. Experimental Details

In the following section, details about the reference implementations of REINFORCE, TRPO and PPO and their parameter settings are summarized for the benchmark experiments and the ablation study. Information about the benchmark environments is given in Sec. B.2

B.1. Algorithm Configurations

The reference implementations of the benchmark algorithms REINFORCE, TRPO and PPO are from the Garage RL framework (Duan et al., 2016). A hyper-parameter grid search has been conducted for each algorithm and each environment on separate random seeds. The parameter ranges and selected hyper-parameters are indicated in Tab. 1. For the benchmark itself, ten runs have been conducted for each algorithm and each environment on the random seeds (404, 931, 159, 380, 858, 708, 16, 448, 136, 989).

The configuration of the DD-OPG method is summarized in Tab. 1.

B.2. Benchmark Environments

The benchmark environments are cartpole, mountaincar and swimmer from the Garage RL framework. Details about the input and state dimensions, as well as the task horizons are listed in Tab. 2.

References

Duan, Y., Chen, X., Houthoofd, R., Schulman, J., and Abbeel, P. Benchmarking deep reinforcement learning

Table 1. Algorithm hyper-parameters for the benchmark tasks.

ALGORITHM	PARAMETER	RANGE	SELECTED
REINFORCE	BATCH SIZE	[400, 5000]	5000
	STEP SIZE	[0.0001, 0.1]	0.03
TRPO	BATCH SIZE	[400, 5000]	5000
	STEP SIZE	[0.0001, 0.1]	0.1
PPO	BATCH SIZE	[400, 5000]	2000
	STEP SIZE	[0.0001, 0.2]	0.2

ALGORITHM	PARAMETER	SYMBOL	SELECTED
DD-OPG	TEMPERATURE	λ	0.1
	PENALTY	γ	0.05
	LENGTHSCALE	$\log \Sigma$	3I
	PATH BUFFER	N_{max}	50

Table 2. Information about the benchmark environments.

ENVIRONMENT	INPUTS	STATES	HORIZON
CARTPOLE	1	4	100
MOUNTAINCAR	1	2	500
SWIMMER	2	13	1000

for continuous control. In *International Conference on Machine Learning (ICML)*, pp. 1329–1338, 2016.

Figure 1. Derivation of the weighted IS policy gradient.

$$\begin{aligned} \nabla_{\theta} \hat{J}^{\text{WIS}}(\theta) &= \nabla_{\theta} \left(\left(\sum_{i=1}^N \frac{p(\tau_i | \theta)}{\frac{1}{N} \sum_j p(\tau_i | \theta_j)} \right)^{-1} \right) \sum_{i=0}^N \frac{p(\tau_i | \theta)}{\frac{1}{N} \sum_j p(\tau_i | \theta_j)} R(\tau_i) + \\ &\quad \left(\sum_{i=1}^N \frac{p(\tau_i | \theta)}{\frac{1}{N} \sum_j p(\tau_i | \theta_j)} \right)^{-1} \sum_{i=0}^N \nabla_{\theta} \left(\frac{p(\tau_i | \theta)}{\frac{1}{N} \sum_j p(\tau_i | \theta_j)} R(\tau_i) \right) \end{aligned} \quad (2)$$

$$\begin{aligned} &= - \left(\sum_{i=1}^N w_i(\theta) \right)^{-2} \left(\sum_{i=1}^N \nabla_{\theta} w_i(\theta) \right) \left(\sum_{i=1}^N w_i(\theta) R(\tau_i) \right) + \\ &\quad \left(\sum_{i=1}^N w_i(\theta) \right)^{-1} \left(\sum_{i=1}^N \nabla_{\theta} w_i(\theta) R(\tau_i) \right) \end{aligned} \quad (3)$$

$$= -\frac{1}{Z^2} \left(\sum_{i=1}^N \nabla_{\theta} w_i(\theta) \right) \left(\sum_{i=1}^N w_i(\theta) R(\tau_i) \right) + \frac{1}{Z} \left(\sum_{i=1}^N \nabla_{\theta} w_i(\theta) R(\tau_i) \right) \quad (4)$$

$$= \frac{1}{Z} \left(\sum_{i=1}^N \nabla_{\theta} w_i(\theta) R(\tau_i) - \sum_{i=1}^N \nabla_{\theta} w_i(\theta) \frac{\sum_{i=1}^N w_i(\theta) R(\tau_i)}{Z} \right) \quad (5)$$

$$= \frac{1}{Z} \left(\sum_{i=1}^N \nabla_{\theta} w_i(\theta) R(\tau_i) - \sum_{i=1}^N \nabla_{\theta} w_i(\theta) \hat{J}^{\text{WIS}}(\theta) \right) \quad (6)$$

$$\nabla_{\theta} \hat{J}^{\text{WIS}}(\theta) = \frac{1}{Z} \sum_{i=1}^N \nabla_{\theta} w_i(\theta) \left(R(\tau_i) - \hat{J}^{\text{WIS}}(\theta) \right) \quad (7)$$

055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109