

---

# Efficient Policy Evaluation with Offline Data Informed Behavior Policy Design

---

Shuze Liu<sup>1</sup> Shangtong Zhang<sup>1</sup>

## Abstract

Most reinforcement learning practitioners evaluate their policies with online Monte Carlo estimators for either hyperparameter tuning or testing different algorithmic design choices, where the policy is repeatedly executed in the environment to get the average outcome. Such massive interactions with the environment are prohibitive in many scenarios. In this paper, we propose novel methods that improve the data efficiency of online Monte Carlo estimators while maintaining their unbiasedness. We first propose a tailored closed-form behavior policy that provably reduces the variance of an online Monte Carlo estimator. We then design efficient algorithms to learn this closed-form behavior policy from previously collected offline data. Theoretical analysis is provided to characterize how the behavior policy learning error affects the amount of reduced variance. Compared with previous works, our method achieves better empirical performance in a broader set of environments, with fewer requirements for offline data.

## 1. Introduction

Reinforcement Learning (RL, Sutton & Barto (2018)) has recently demonstrated great success in solving sequential decision-making problems. For example, AlphaStar (Vinyals et al., 2019) defeats the best human StarCraft II players and is ranked at the GrandMaster level in the StarCraft ladder. The canonical RL paradigm behind the success, however, requires massive active interactions with the environment to obtain data (Sutton, 1988; Watkins & Dayan, 1992; Sutton et al., 1999; Mnih et al., 2015). Those data are called online data, and this paradigm is called online RL. Requiring massive online data is, however, prohibitive in many scenarios. First, obtaining massive online data can

be both expensive and slow in the real world (Li, 2019; Zhang, 2023). Second, even if a simulator is available, obtaining massive online data can still be prohibitively slow for high-fidelity simulation (Chervonyi et al., 2022).

Offline RL (Ernst et al., 2005; Lange et al., 2012; Fujimoto et al., 2019; Levine et al., 2020) attacks this issue using existing, previously logged data, called offline data. Compared with online data, offline data is cheaper and safer (Li, 2019; Zhang, 2023). Offline RL also demonstrates great success. For example, Mathieu et al. (2023) train an offline AlphaStar, which uses only existing human replays without any interaction with the StarCraft II simulator during training. The offline AlphaStar obtains over 90% win rates against the supervised learning agent in Vinyals et al. (2019).

However, most RL practitioners, even offline RL practitioners, still heavily rely on online Monte Carlo estimators. For example, Mathieu et al. (2023) repeatedly execute their trained offline AlphaStar agents in the StarCraft II simulator and use the win rates as the performance metric for hyperparameter tuning and evaluating different algorithmic design choices. This evaluation practice is the straightforward online Monte Carlo evaluation and requires massive online data. There are indeed offline evaluation methods, most of which, however, still rely on online Monte Carlo evaluation for hyperparameter tuning and testing different algorithmic design choices (see, e.g., Fu et al. (2020); Gülçehre et al. (2020); Schrittwieser et al. (2021); Mathieu et al. (2023)).

**Improving the sample efficiency of online Monte Carlo estimators while maintaining their unbiasedness** is thus a need for both online and offline RL practitioners. We emphasize *unbiasedness* because it is arguably one of the key reasons that make Monte Carlo estimators so dominating. In this paper, we make three contributions toward fulfilling this need. **First**, we propose tailored closed-form behavior policies that *provably* reduce the variance of online Monte Carlo estimators. **Second**, we design efficient algorithms to learn the closed-form behavior policies from offline data. Theoretical analysis is provided to characterize how the behavior policy learning error affects the amount of reduced variance. Notably, this learning error does not introduce any bias in the estimation. **Third**, we conduct thorough empirical studies in a broad set of environments. Compared with previous works, our method achieves better empirical

---

<sup>1</sup>Department of Computer Science, University of Virginia. Correspondence to: Shuze Liu <shuzeliu@virginia.edu>.

performance while being less restrictive on offline data.

## 2. Background

We consider a finite horizon Markov Decision Process (MDP, [Puterman \(2014\)](#)) with a finite state space  $S$ , a finite action space  $A$ , a reward function  $r : S \times A \rightarrow \mathbb{R}$ , a transition probability function  $p : S \times S \times A \rightarrow [0; 1]$ , an initial distribution  $p_0 : S \rightarrow [0; 1]$ , and a constant horizon length  $T$ . Without loss of generality, we consider the undiscounted setting for simplifying notations. Our results naturally apply to the discounted setting ([Puterman, 2014](#)) as long as the horizon is fixed and finite. For any integer  $n$ , we define as shorthand  $[n] \doteq \{0; 1; \dots; n\}$ . At time step 0, an initial state  $S_0$  is sampled from  $p_0$ . At time step  $t \in [T - 1]$ , an action  $A_t$  is sampled according to  $\pi_t(\cdot | S_t)$  where  $\pi_t : A \times S \rightarrow [0; 1]$  is the policy at time step  $t$ . A finite reward  $R_{t+1} \doteq r(S_t; A_t)$  is then emitted and a successor state  $S_{t+1}$  is sampled from  $p(\cdot | S_t; A_t)$ . We define abbreviations  $f_{i:j} \doteq \{f_i; \dots; f_j\}$  and  $\beta_{i:T} \doteq \{\beta_i; \dots; \beta_T\}$ . The return at time step  $t$  is defined as  $G_t \doteq \sum_{i=t+1}^T R_i$  which allows defining the state- and action-value functions as  $v_{\pi_t}(s) \doteq \mathbb{E}[G_t | S_t = s]$  and  $q_{\pi_t}(s; a) \doteq \mathbb{E}[G_t | S_t = s; A_t = a]$ . We use the total rewards performance metric ([Puterman, 2014](#)) to measure the performance of the policy  $\pi$ , which is defined as  $J(\pi) \doteq \sum_s p_0(s) v_{\pi}(s)$ . In this paper, we focus on Monte Carlo methods introduced by [Kakutani \(1945\)](#) to estimate the total rewards  $J(\pi)$ . Among its variants, the most straightforward and widely used way is to draw samples of  $J(\pi)$  by executing the policy  $\pi$  online. As the number of samples increases, the empirical average of the sampled returns converges to  $J(\pi)$ . This idea is called on-policy learning ([Sutton 1988](#)) because it estimates a policy  $\pi$  by executing itself.

From now on, we consider off-policy learning, where we estimate the total rewards  $J(\pi)$  of an interested policy  $\pi$ , called the target policy, by executing a different policy  $\pi'$ , called the behavior policy. In off-policy learning, each trajectory  $f_{S_0; A_0; R_1; S_1; A_1; R_2; \dots; S_{T-1}; A_{T-1}; R_T g$  is generated by a behavior policy  $\pi'$  with  $S_0 \sim p_0; A_t \sim \pi'_t(\cdot | S_t); t \in [T - 1]$ . Let

$$f_{t:T-1} \doteq \{f_{S_t; A_t; R_{t+1}; \dots; S_{T-1}; A_{T-1}; R_T g}$$

be a shorthand for a segment of a random trajectory generated by the behavior policy  $\pi'$  from the time step  $t$  to the time step  $T - 1$  inclusively. In off-policy learning, we use the importance sampling ratio to reweight rewards collected by  $\pi'$  in order to give an estimate of  $J(\pi)$ . The importance sampling ratio at time step  $t$  is defined as  $\rho_t \doteq \frac{\pi_t(A_t | S_t)}{\pi'_t(A_t | S_t)}$ . The product of importance sampling ratios from time  $t$  to  $t^0$  is defined as  $\rho_{t:t^0} \doteq \prod_{k=t^0}^t \frac{\pi_k(A_k | S_k)}{\pi'_k(A_k | S_k)}$ . There are various ways to use the importance sampling ratios in off-policy

learning ([Geweke, 1988](#); [Hesterberg, 1995](#); [Koller & Friedman, 2009](#); [Thomas, 2015](#)). We start with the per-decision importance sampling estimator (PDIS, [Precup et al. \(2000\)](#)) in this work and leave the investigation of others for future work. The PDIS Monte Carlo estimator is defined as

$$G^{\text{PDIS}}(f_{t:T-1}) \doteq \prod_{k=t}^{T-1} \rho_k R_{k+1} \quad (1)$$

and is unbiased for any behavior policy  $\pi'$  that covers target policy  $\pi$  ([Precup et al., 2000](#)). In other words, when  $\beta_s, \beta_a, \pi'_t(a|s) = 0 \Rightarrow \pi_t(a|s) = 0$ , we have  $\beta_s, \beta_a$ ,

$$\mathbb{E}[G^{\text{PDIS}}(f_{t:T-1}) | S_t = s] = v_{\pi_t}(s)$$

We intensively use the recursive form of the PDIS estimator:

$$G^{\text{PDIS}}(f_{t:T-1}) = \begin{cases} \rho_t R_{t+1} + G^{\text{PDIS}}(f_{t+1:T-1}) & t \in [T - 2]; \\ \rho_t R_{t+1} & t = T - 1; \end{cases} \quad (2)$$

Since the PDIS estimator is unbiased, reducing its variance is sufficient for improving its sample efficiency. We achieve this variance reduction by designing and learning proper behavior policies.

## 3. Variance Reduction in Statistics

In this section, we provide the mathematical foundation for variance reduction with importance sampling ratios. The notations here are independent of the rest of this paper. We use similar notations only for easy interpretation in later sections. Consider a discrete random variable  $A$  taking values from a finite space  $A$  according to a probability mass function  $\pi : A \rightarrow [0; 1]$  and a function  $q : A \rightarrow \mathbb{R}$  mapping a value in  $A$  to a real number. We are interested in estimating  $\mathbb{E}_A[q(A)]$ . The ordinary Monte Carlo methods then sample  $f_{A_1; \dots; A_N g}$  from  $\pi$  and use the empirical average  $\frac{1}{N} \sum_{i=1}^N q(A_i)$  as the estimate. In statistics, importance sampling is introduced as a variance reduction technique for Monte Carlo methods ([Rubinstein 1981](#)). The main idea is to sample  $f_{A_i; \dots; A_N g}$  from a different distribution  $\pi'$  and use  $\frac{1}{N} \sum_{i=1}^N \rho(A_i) q(A_i)$  as the estimate, where  $\rho(A) \doteq \frac{\pi(A)}{\pi'(A)}$  is the importance sampling ratio. Assuming  $\pi'$  covers  $\pi$ , i.e.,

$$\beta_a; \pi'_t(a) = 0 \Rightarrow \pi_t(a) = 0; \quad (3)$$

the importance sampling ratio weighted empirical average is then unbiased, i.e.,

$$\mathbb{E}_A[q(A)] = \mathbb{E}_A[\rho(A)q(A)];$$

If the sampling distribution  $\pi'$  is carefully designed, the variance can also be reduced. To adapt this idea for RL, we relax the condition (3) in this section. We formulate

this problem of searching a variance-reducing sampling distribution as an optimization problem:

$$\min_{\pi} \text{Var}_{\pi} \left( \sum_{A \in \mathcal{A}} q(A) \right) \quad (4)$$

Here  $\pi$  denotes the set of all the policies that give unbiased estimations, i.e.,

$$\pi \doteq \left\{ \pi \mid \mathbb{E}_{\pi} \left[ \sum_{A \in \mathcal{A}} q(A) \right] = \mathbb{E}_{\pi} [q(A)] \right\};$$

where  $\mathcal{X}$  denotes the set of all probability distributions on the set  $\mathcal{X}$ . Solving (4) is actually very challenging. To see this, consider a concrete example where  $\mathcal{A} = \{a_1, a_2, a_3\}$  and

$$\begin{aligned} \mathcal{A} &= \{a_1, a_2, a_3\} \\ q(a_1) &= 10 & q(a_1) &= 0:1 & q(a_1) &= 0 \\ q(a_2) &= 2 & q(a_2) &= 0:5 & q(a_2) &= 0 \\ q(a_3) &= 2 & q(a_3) &= 0:4 & q(a_3) &= 1 \end{aligned} \quad (5)$$

It can be computed that  $\mathbb{E}_{\pi} [q(A)] = 0:8$  and  $\mathbb{E}_{\pi} \left[ \sum_{A \in \mathcal{A}} q(A) \right] = 0:8$ . In other words, we could sample  $A$  from  $\pi$  and use  $\sum_{A \in \mathcal{A}} q(A)$  as an estimator. This estimator is unbiased. But apparently, this does not cover. Moreover, since  $\pi$  is deterministic, the variance of this estimator is 0. Then  $\pi$  is an optimal sampling distribution. However,  $\pi$  is hand-crafted based on the knowledge that  $q(a_1) + q(a_2) + q(a_3) = 0$ . Without such knowledge, we argue that there is little hope to find this. This example suggests that searching over the entire might be too ambitious. One natural choice presented by Rubinstein (1981) is to restrict the search to

$$\pi \doteq \left\{ \pi \mid \sum_{A \in \mathcal{A}} q(A) = 0 \Rightarrow \sum_{A \in \mathcal{A}} q(A) = 0 \right\} \quad (6)$$

In other words, we aim to find a variance-minimizing sampling distribution among all distributions that cover. Because coverage implies unbiasedness, we have  $\pi \subseteq \pi$ . In this work, we enlarge  $\pi$  to be defined as

$$\pi \doteq \left\{ \pi \mid \sum_{A \in \mathcal{A}} q(A) = 0 \Rightarrow \sum_{A \in \mathcal{A}} q(A) = 0 \right\} \quad (7)$$

following Owen (2013). The space weakens the assumption in (6). Owen (2013) proves that any distribution  $\pi$  gives unbiased estimation, though may not cover.

Lemma 1.  $\mathcal{A} = \{a_1, a_2, a_3\}$ ;  $\mathbb{E}_{\pi} \left[ \sum_{A \in \mathcal{A}} q(A) \right] = \mathbb{E}_{\pi} [q(A)]$ :

For completeness, its proof is in Appendix A.1. We now consider the variance minimization problem on i.e.,

$$\min_{\pi} \text{Var}_{\pi} \left( \sum_{A \in \mathcal{A}} q(A) \right) \quad (8)$$

The following lemma from Owen (2013) gives an optimal solution to the optimization problem (8).

Lemma 2. Define  $\pi(a) / \pi(a)q(a)$ . Then  $\pi$  is an optimal solution to (8).

For completeness, its proof is detailed in Appendix A.2. Here by

$$\pi(a) / \pi(a)q(a)$$

with some non-negative  $w(a)$ , we mean

$$\pi(a) \doteq \frac{\pi(a)w(a)}{\sum_b \pi(b)w(b)}$$

The reader may notice that if  $\pi(a)w(a) = 0$  for all  $a$ , the above ‘‘reweighted’’ distribution is not well defined. We then use the convention to interpret  $\pi(a)$  as a uniform distribution, i.e.,  $\pi(a) = 1/|\mathcal{A}|$ . We adopt this convention in using  $\pi$  in the rest of the paper to simplify the presentation. The following lemma gives intuition on the optimality of  $\pi$ , whose proof is in Appendix A.3.

Lemma 3. If  $\sum_{A \in \mathcal{A}} q(A) = 0$  or  $\sum_{A \in \mathcal{A}} q(A) > 0$ , then  $\pi$  is optimal, and the one defined in Lemma 2 gives a zero variance, i.e.,  $\text{Var}_{\pi} \left( \sum_{A \in \mathcal{A}} q(A) \right) = 0$ .

An optimal sampling distribution proportional to  $\pi(a)q(a)$  dates back to Kahn & Marshall (1953); Rubinstein (1981); Benjamin Melamed (1998) and is commonly used in RL (Carpentier et al., 2015; Mukherjee et al., 2022). We, however, make two remarks. First, we show such a sampling distribution can be suboptimal in  $\pi$ . For (5), such a sampling distribution incurs strictly positive variance, but in (5) has a zero variance and is also unbiased. Second, different from existing literature in RL (Carpentier et al., 2015; Sutton & Barto, 2018; Mukherjee et al., 2022), our definition in Lemma 2 does not need to cover. Nevertheless, we note that Lemma 1 still ensures that gives unbiased estimation (Owen, 2013) and extend unbiasedness to RL settings in Theorem 1.

#### 4. Variance Reduction in Reinforcement Learning

We now apply the techniques in Section 3 in RL. In particular, we seek to reduce the variance of  $G^{\text{PDIS}} \left( \begin{smallmatrix} 0 & T & 1 \\ 0 & T & 1 \end{smallmatrix} \right)$  by designing a proper behavior policy. Of course, we need to ensure that the PDIS estimator with this behavior policy is unbiased. In other words, ideally we should search over

$$\pi \doteq \left\{ \pi \mid \mathbb{E}_{\pi} \left[ G^{\text{PDIS}} \left( \begin{smallmatrix} 0 & T & 1 \\ 0 & T & 1 \end{smallmatrix} \right) \right] = J(\pi) \right\}$$

As discussed in Section 3, this is too ambitious without domain-specific knowledge. Instead, we can search over all policies that cover, i.e.,

$$\pi \doteq \left\{ \pi \mid \sum_{s \in \mathcal{S}} \pi(s)q(s) = 0 \Rightarrow \sum_{s \in \mathcal{S}} \pi(s)q(s) = 0 \right\}$$

The set  $\pi$  contains all policies that satisfy the policy coverage constraint in off-policy learning (Sutton & Barto 2018). Similar to (7), we can also enlarge  $\pi$  to

$$\pi \doteq \left\{ \pi \mid \sum_{s \in \mathcal{S}} \pi(s)q(s) = 0 \Rightarrow \sum_{s \in \mathcal{S}} \pi(s)q(s) = 0 \right\}$$

The following theorem ensures the desired unbiasedness, its proof is in Appendix A.5. We are now ready to define  $\pi_t$  which is proved in Appendix A.4.

**Theorem 1 (Unbiasedness)**  $\forall t \in \{0, \dots, T-1\}$ ,  $\forall s \in \mathcal{S}$ ,  $\mathbb{E} G^{\text{PDIS}}( \pi_t ) | \mathcal{S}_t = s = v_t(s)$ .

One immediate consequence of Theorem 1 is that  $\pi_t$  is unbiased. In this paper, we consider a set  $\Pi_t$  such that  $\pi_t \in \Pi_t$ . This inherits the unbiasedness property of  $\pi_t$  and is less restrictive than  $\Pi_t$ , the classical search space of behavior policies. This will be defined shortly. We now formulate our problem as

$$\min_{\pi_t} \mathbb{V} G^{\text{PDIS}}( \pi_t ) : \quad (9)$$

By the law of total variance, for any  $\pi_t$ , we decompose the variance of the PDIS estimator as

$$\begin{aligned} & \mathbb{V} G^{\text{PDIS}}( \pi_t ) \\ &= \mathbb{E}_{S_0} \mathbb{V} G^{\text{PDIS}}( \pi_t ) | S_0 \\ & \quad + \mathbb{V}_{S_0} \mathbb{E} G^{\text{PDIS}}( \pi_t ) | S_0 \\ &= \mathbb{E}_{S_0} \mathbb{V} G^{\text{PDIS}}( \pi_t ) | S_0 + \mathbb{V}_{S_0} (v_t(S_0)) : \\ & \quad \text{(by Theorem 1)} \end{aligned}$$

The second term  $\mathbb{V}_{S_0} (v_t(S_0))$  is a constant given a target policy  $\pi_t$  and is unrelated to the choice of  $\pi_t$ . In the first term, the expectation is taken over  $S_0$  that is determined by the initial probability distribution  $p_0$ . Consequently, to solve the problem (9), it is sufficient to solve for each  $s$

$$\min_{\pi_t} \mathbb{V} G^{\text{PDIS}}( \pi_t ) | S_0 = s : \quad (10)$$

Denote the variance of the state value for the next state given the current state-action pair  $(s; a)$  as  $\sigma_{t+1}^2(s; a)$ . We have  $\sigma_{t+1}^2(s; a) = 0$  for  $t = T-1$  and otherwise

$$\sigma_{t+1}^2(s; a) \doteq \mathbb{V}_{S_{t+1}} (v_{t+1}(S_{t+1}) | S_t = s, A_t = a) : \quad (11)$$

We now construct a behavior policy  $\pi_t$  as

$$\pi_t(a|s) \propto \pi_t(a|s) \mathbb{P} \overline{u_{t+1}(s; a)} : \quad (12)$$

where  $\overline{u_{t+1}(s; a)} \doteq \sigma_{t+1}^2(s; a)$  for  $t = T-1$  and otherwise

$$\begin{aligned} \overline{u_{t+1}(s; a)} &= \sigma_{t+1}^2(s; a) + \sigma_{t+1}^2(s; a) \\ & \quad + \mathbb{P}_{s^0} \mathbb{V} G^{\text{PDIS}}( \pi_{t+1} ) | S_{t+1} = s^0 : \end{aligned} \quad (13)$$

Notably,  $\pi_t$  and  $\overline{u_{t+1}}$  are defined backwards and alternatively, i.e., they are defined in the order of  $u_{T-1}; \pi_{T-1}; u_{T-2}; \pi_{T-2}; \dots; u_0; \pi_0$ . We prove  $\pi_t$  is optimal in the following sense.

**Theorem 2 (Optimal Behavior Policy)** For any  $\pi_t$  and  $s$ , the behavior policy  $\pi_t(a|s)$  defined above is an optimal solution to the following problem

$$\min_{\pi_t} \mathbb{V} G^{\text{PDIS}}( \pi_t ) | S_t = s ;$$

where  $\pi_t \doteq \pi_{t-2} \in \Pi_{t-2} \cap \Pi_{t-1}$  and  $\pi_t(a|s) = 0 \Rightarrow \pi_t(a|s) \overline{u_{t+1}(s; a)} = 0$ .

Theorem 2 indicates that  $\pi_t$  achieves optimality for the optimization problem (10). Since  $\overline{u_{t+1}(s; a)} = 0 \Rightarrow \pi_t(a|s) = 0$  by the non-negativity of the summands in (13), we have  $\pi_t(a|s) = 0 \Rightarrow \pi_t(a|s) \overline{u_{t+1}(s; a)} = 0$ . This indicates that the set of policies considered in Theorem 2 are unbiased and includes at least all the policies that cover the target policy, which is the classical behavior policy search space (Precup et al., 2000; Maei, 2011; Sutton et al., 2016; Zhang, 2022).

Unfortunately, empirically implementing  $\pi_t$  requires knowledge of  $\overline{u_{t+1}}$  (13) that contains the transition function. Approximating the transition function is very challenging in MDPs with large stochasticity and function approximation (cf. model-based RL (Sutton, 1990; Sutton et al., 2008; Deisenroth & Rasmussen, 2011; Chua et al., 2018)). Thus, we seek to build another policy that can be easily implemented without direct knowledge of the transition function  $p$  (cf. model-free RL (Sutton, 1988; Watkins, 1989)).

We achieve this by aiming at one-step optimality instead of global optimality. We try to find the best  $\pi_t$  assuming in the future we follow  $\pi_{t+1}; \dots; \pi_{T-1}$ , instead of  $\pi_{t+1}; \dots; \pi_{T-1}$ . We refer to this one-step optimal behavior policy as  $\hat{\pi}_t$ . Similarly, to define optimality, we first need to specify the set of policies we are concerned about. To this end, we define

$$\hat{q}_{t+1}(s; a) \doteq \sigma_{t+1}^2(s; a) \quad (14)$$

for  $t = T-1$  and otherwise

$$\begin{aligned} \hat{q}_{t+1}(s; a) &\doteq \sigma_{t+1}^2(s; a) + \hat{q}_{t+1}(s; a) \\ & \quad + \mathbb{P}_{s^0} \mathbb{V} G^{\text{PDIS}}( \pi_{t+1} ) | S_{t+1} = s^0 : \end{aligned} \quad (15)$$

Notably,  $\hat{q}_{t+1}(s; a)$  is always non-negative since all the summands are non-negative. Accordingly, we define for  $t \in [T-1]$ ,  $\hat{\pi}_t \doteq \pi_{t-2} \in \Pi_{t-2} \cap \Pi_{t-1}$  and  $\hat{\pi}_t(a|s) = 0 \Rightarrow \hat{\pi}_t(a|s) \hat{q}_{t+1}(s; a) = 0$ . Comparing (13) and (15), the optimality of  $\pi_t$  implies that for  $s; a; t$ , we have  $\pi_t(a|s) \overline{u_{t+1}(s; a)} = 0$ . As a result, if  $\pi_t \in \hat{\Pi}_t$ , we have

$$\begin{aligned} \pi_t(a|s) = 0 &\Rightarrow \pi_t(a|s) \hat{q}_{t+1}(s; a) = 0 \\ &\Rightarrow \pi_t(a|s) \overline{u_{t+1}(s; a)} = 0 ; \end{aligned}$$

indicating  $\pi_t \in \hat{\Pi}_t$ . In other words, we have  $\hat{\Pi}_t \subseteq \Pi_t$ . To search for  $\pi_{0:T-1}$ , we work on  $\hat{\Pi} \doteq \hat{\Pi}_0 \times \dots \times \hat{\Pi}_{T-1}$ . To summarize, we have  $\hat{\Pi} \subseteq \Pi$ . Recall that  $\hat{\Pi}$  is the set of all behavior policies such that the corresponding PDIS estimator is unbiased as a sufficient but not necessary condition to ensure such unbiasedness (Theorem 1).  $\hat{\Pi}$  is a restriction of  $\Pi$  such that we are able to find an optimal solution. We restrict  $\Pi$  to  $\hat{\Pi}$ , aiming for a

sub-optimal but implementable policy), is still larger than  $V_{t, \pi_t}^{\text{PDIS}}(f_{t:T-1}^{t+1:T-1})$ , which is the space with the coverage assumption that previous works (Precup et al., 2000; Maei, 2011; Sutton et al., 2016; Sutton & Barto, 2018; Zhang, 2022) consider. After confirming the space of behavior policies, we formulate the optimization problem for designing an efficient behavior policy to achieve one-step optimality as

$$\min_{\pi_t \in \Pi_t} V_{t, \pi_t}^{\text{PDIS}}(f_{t:T-1}^{t+1:T-1}) \mid S_t = s_t \quad (16)$$

According to the recursive expression of the variance in Lemma 4 in Appendix A.5, we rewrite (16) as

$$\min_{\pi_t \in \Pi_t} E_{A_t \sim \pi_t} \left[ E_{S_{t+1}} \left[ V_{t+1, \pi_{t+1}}^{\text{PDIS}}(f_{t+1:T-1}^{t+1:T-1}) \mid S_{t+1} \right] \mid S_t; A_t \right] + \sigma_{t+1}^2(S_t; A_t) \mid S_t; \quad (17)$$

where the objective can be further simplified as

$$\begin{aligned} & E_{A_t \sim \pi_t} \left[ E_{S_{t+1}} \left[ V_{t+1, \pi_{t+1}}^{\text{PDIS}}(f_{t+1:T-1}^{t+1:T-1}) \mid S_{t+1} \right] \mid S_t; A_t \right] + \sigma_{t+1}^2(S_t; A_t) \mid S_t \\ &= E_{A_t \sim \pi_t} \left[ \frac{2}{q} \hat{q}_{t+1}(S_t; A_t) \mid S_t \right] \quad (\text{By (15)}) \\ &= V_{A_t \sim \pi_t} \left[ \frac{2}{q} \hat{q}_{t+1}(S_t; A_t) \mid S_t \right] \\ &= E_{A_t \sim \pi_t}^2 \left[ \frac{q}{\hat{q}_{t+1}(S_t; A_t) \mid S_t} \right] \quad (\text{Lemma 1 and } t \geq t) \end{aligned}$$

Since the second term is unrelated to  $\pi_t$  it is equivalent to solving

$$\min_{\pi_t \in \Pi_t} V_{A_t \sim \pi_t} \left[ \frac{q}{\hat{q}_{t+1}(S_t; A_t) \mid S_t} \right]$$

According to Lemma 2,

$$\hat{\pi}_t(a \mid s) / \pi_t(a \mid s) \propto \frac{q}{\hat{q}_{t+1}(s; a)} \quad (18)$$

is an optimal solution to (17). We now present our main result that  $\hat{\pi}_t$  provably reduces variance.

**Theorem 3 (Variance Reduction)** For any  $\pi_t$  and  $\hat{\pi}_t$ ,

$$\begin{aligned} & V_{t, \pi_t}^{\text{PDIS}}(f_{t:T-1}^{t+1:T-1}) \mid S_t = s_t \\ & V_{t, \hat{\pi}_t}^{\text{PDIS}}(f_{t:T-1}^{t+1:T-1}) \mid S_t = s_t \leq \pi_t(s_t): \end{aligned}$$

To define  $\pi_t(s)$ , first define  $\pi_t(s) =$

$$\prod_{a \sim \pi_t(a \mid s)} \hat{q}_{t+1}(s; a) \propto \prod_{a \sim \pi_t(a \mid s)} \frac{q}{\hat{q}_{t+1}(s; a)}^2$$

Then we define  $\pi_t(s) \propto \pi_t(s)$  for  $t = T-1$  and otherwise

$$\pi_t(s) \propto \pi_t(s) + E_{A_t \sim \pi_t} \left[ E_{S_{t+1}} \left[ \pi_{t+1}(S_{t+1}) \mid S_t; A_t \right] \right] \quad (19)$$

Its proof is in Appendix A.6. Notably, this is always non-negative by Jensen's inequality, ensuring the non-negativity of  $\pi_t$  and thus the variance reduction property. Moreover,  $\pi_t(s) = 0$  occurs only when all actions have the same  $\hat{q}_{t+1}$  on the states. It is reasonable to conjecture that this is rare in practice. So  $\pi_t(s)$  is likely to be strictly positive. This shows the variance of the PDIS estimator with a state  $s$  is provably smaller than or equal to that with the straightforward on-policy Monte Carlo estimator, by at least  $\pi_t(s)$ . The magnitude of  $\pi_t(s)$  depends on a specific target policy and the environment. We empirically show the variance reduction is significant in commonly used benchmarks in Section 7.

## 5. Learning Closed-Form Behavior Policies

We now present efficient algorithms to learn the closed-form behavior policy  $\hat{\pi}_t$ . Despite that  $\hat{q}_{t+1}$  in (15) has a complicated definition, we prove that it has a concise representation. It is exactly the action value function of the policy with the same transition function but a different reward function.

**Theorem 4.** Define

$$\hat{r}_{t+1}(s; a) \doteq 2r(s; a) \hat{q}_{t+1}(s; a) - r^2(s; a) \quad (20)$$

Then  $\hat{q}_{t+1}(s; a) = \hat{r}_{t+1}(s; a)$  for  $t = T-1$  and otherwise

$$\begin{aligned} & \hat{q}_{t+1}(s; a) = \sum_{s^0; a^0} p(s^0 \mid s; a) \pi_{t+1}(a^0 \mid s^0) \hat{q}_{t+1}(s^0; a^0) \\ &= \hat{r}_{t+1}(s; a) + \sum_{s^0; a^0} p(s^0 \mid s; a) \pi_{t+1}(a^0 \mid s^0) \hat{q}_{t+1}(s^0; a^0) \end{aligned} \quad (21)$$

Its proof is in Appendix A.7. This observation makes it possible to apply any off-the-shelf of line policy evaluation methods to learn  $\hat{q}_{t+1}$ , after which the behavior policy can be computed easily with (18). For generality, we consider the behavior policy agnostic of line learning setting (Nachum et al., 2019), where the offline data in the form of  $\{(t_i; s_i; a_i; r_i; s_i^0)\}_{i=1}^n$  consists of previously logged data tuples. In the  $i$ -th data tuple,  $t_i$  is the time step,  $s_i$  is the state at time step  $t_i$ ,  $a_i$  is the action executed on state  $s_i$ ,  $r_i$  is the sampled reward, and  $s_i^0$  is the successor state. Those tuples can be generated by one or more, known or unknown behavior policies. Those tuples do not need to form a complete trajectory.

In this paper, we choose Fitted Evaluation (FQE, Le et al. (2019)) as a demonstration, but our framework is ready to incorporate any state-of-the-art of line policy evaluation methods to approximate  $\hat{q}_{t+1}$ . To learn  $\hat{\pi}_t$ , it is sufficient to learn  $\hat{r}_{t+1}$  and  $\hat{q}_{t+1}$ . FQE can be used to learn  $\hat{r}_{t+1}$  and learning is a simple regression problem. FQE is then invoked again w.r.t. the learned  $\hat{r}_{t+1}$  to learn an approximation  $\hat{q}_{t+1}$ . We refer the reader to Algorithm 1 for a detailed exposition of our algorithm. We split the offline data into training

**Algorithm 1** Offline Data Informed (ODI) algorithm

- 1: Input: Estimators  $\hat{r}(s; a)$ ,  $q_{t,t}(s; a)$ ,  $\hat{q}_{t,t}(s; a)$ , a target policy  $\pi$ , an offline dataset  $D = \{f(t_i; s_i; a_i; r_i; s_i)\}_{i=1}^m$
- 2: Output: a behavior policy  $\hat{\pi}$
- 3: Approximate  $r$  from  $D$  using supervised learning
- 4: Approximate  $q_{t,t}$  from  $D$  using any of the RL method (e.g. Fitted Q-Evaluation)
- 5: Compute  $\hat{r}_t$  by (20) for each data pair  $i \in D$
- 6: Construct  $D_{\hat{r}} = \{f(t_i; s_i; a_i; \hat{r}_i; s_i)\}_{i=1}^m$  by plugging  $\hat{r}_t$  into  $D$
- 7: Approximate  $\hat{q}_{t,t}$  from  $D_{\hat{r}}$  by (21) using any of the RL method (e.g. Fitted Q-Evaluation)
- 8: Return:  $\hat{\pi}_t(a|s) / \pi_t(a|s)$ ,  $\hat{q}_{t,t}(s; a)$

sets and test sets to tune all the hyperparameters of in Algorithm 1, based on the supervised learning loss or the FQE loss on the test set. We remark that FQE loss on the test set is known to be an inaccurate signal (Fujimoto et al., 2022) so our estimation would be poorly tuned in this sense. We, however, notice that even with such a poorly tuned estimation, the variance reduction in the tested environments is still significant. This suggests that  $\hat{q}_{t,t}(s; a)$  in Theorem 3 is likely to be large and demonstrates the robustness of our approach. Since  $\hat{q}_{t,t}(s; a)$  is proved to be always non-negative (cf. (15)), we use positive function class for FQE in approximating  $q$ , e.g., a neural network with softplus as the last activation function.

In the following, we theoretically analyze how the error in approximating  $q$  affects the amount of reduced variance in Theorem 3. We assume  $\hat{q}_{t,t}(s; a)$  is not only non-negative but also positive. Given its non-negative summand (15), we argue that this positivity assumption is not restrictive at all. We use  $\hat{q}_{t,t}^+(s; a) > 0$  to denote our approximation to  $\hat{q}_{t,t}(s; a)$ . The approximation error can then be captured by

$$\hat{q}_{t,t}(s; a) \doteq \hat{q}_{t,t}^+(s; a) - \epsilon_{t,t}(s; a) > 0: \quad (22)$$

If  $\pi_t(s; a)$  is 1, there is no approximation error (cf. (22)). The actual learned behavior policy is then denoted by

$$\hat{\pi}_t^+(a|s) / \pi_t(a|s) = \frac{\hat{q}_{t,t}^+(s; a)}{\hat{q}_{t,t}(s; a)}: \quad (23)$$

Then, we generalize Theorem 3 to the following theorem.

**Theorem 5.** For any  $\pi$  and  $\pi^+$ ,

$$\begin{aligned} V(G_{t:T}^{\text{PDIS}}(\hat{\pi}_{t:T}^+ | \pi^+) | S_t = s) \\ V(G_{t:T}^{\text{PDIS}}(\pi | \pi^+) | S_t = s) - \hat{c}_t^+(s): \end{aligned}$$

To denote  $\hat{c}_t^+(s)$ , first denote

$$\hat{c}_t^+(s) \doteq \mathbb{P}_{a \sim \pi_t(a|s)} \hat{q}_{t,t}(s; a)$$

$$\begin{aligned} \mathbb{P}_{a \sim \pi_t(a|S_t)} \mathbb{P}_{\pi^+} \frac{\hat{q}_{t,t}(S_t; a)}{\pi^+(S_t; a)} \\ \mathbb{P}_{a \sim \pi_t(a|S_t)} \mathbb{P}_{\pi^+} \frac{1}{\pi^+(S_t; a)} \mathbb{P}_{\pi^+} \hat{q}_{t,t}(S_t; a): \end{aligned}$$

Then we denote  $\hat{c}_t^+(s) \doteq c_t^+(s)$  for  $t = T - 1$  and otherwise

$$\begin{aligned} \hat{c}_t^+(s) \\ \doteq c_t^+(s) + E_{A_t \sim \hat{\pi}_t^+} [E_{S_{t+1}} [c_{t+1}^+(S_{t+1}) | s; A_t]]: \end{aligned} \quad (24)$$

Its proof is in Appendix A.8. When there is no estimation error, i.e.,  $\hat{q}_{t,t}(s; a) = q_{t,t}(s; a)$ ,  $\hat{c}_t^+$  and  $c_t^+$  reduce to  $c_t$  and  $c_t$  in Theorem 3, which is non-negative by Jensen's inequality. As discussed earlier, it is reasonable to conjecture that  $\hat{c}_t^+(s)$  is likely to be strictly positive. This leaves room to tolerate estimation errors such that  $\hat{c}_t^+(s)$  can still be positive even if  $\hat{q}_{t,t}(s; a) \notin 1$ . Because the sign of  $\hat{c}_t^+$  only depends on the current  $\hat{q}_{t,t}$ , the estimation error in the future step does not affect current  $\hat{c}_t^+$ . Notably, even if some  $c_{t+1}^+(S_{t+1}) < 0$ ,  $\hat{c}_t^+(s)$  can still be positive. This is because  $\hat{c}_t^+(s)$  depends on the expectation of the  $\hat{c}_{t+1}^+(S_{t+1})$ , not a single value, and  $\hat{c}_t^+$  can still be positive. This makes our approach robust to the approximation error. It is important to note that the PDIS estimator with  $\hat{r}_t(a|s)$  is always unbiased, regardless of the approximation error.

Theorem 5 makes it straightforward to analyze how the offline data affects the amount of the reduced variance. For example, if FQE is used, one can resort to Munos (2003); Antos et al. (2008); Munos & Szepesvári (2008); Chen & Jiang (2019) to connect offline data and the approximation error. Theorem 5 then directly relates the approximation error to the amount of reduced variance. We, however, omit such analysis since it deviates from our main contribution.

## 6. Related Work

Monte Carlo methods. Reducing the variance of Monte Carlo estimators via learning a proper behavior policy has been explored before. Hanna et al. (2017) model the problem of finding a variance-reducing behavior policy as an optimization problem and thus rely on stochastic gradient descent to update a parameterized behavior policy. In particular, Hanna et al. (2017) consider the ordinary importance sampling. By contrast, we consider the per-decision importance sampling, which is fundamentally better (Precup et al., 2000). Moreover, Hanna et al. (2017) require new online data to learn this behavior policy. By contrast, our method works with offline data and does not need any online data for behavior policy learning. Hanna et al. (2017) also require the online data to be complete trajectory. By contrast, our method copes well with incomplete offline tuples. Sukherjee et al. (2022) also investigate variance-reducing behavior policies for the per-decision importance sampling estimator. Their results, however, apply to only tree-structured MDPs,

	MDP	Data to learn	Parameterization of	Gridworld size	Other environments
Ours	general	offline data	no assumption	27,000	MuJoCo robotics
BPS (Hanna et al., 2017)	general	online data	need to be known	1,600	CartPole, Acrobot
ROS (Zhong et al., 2022)	general	online data	need to be known	1,600	CartPole
ReVar (Mukherjee et al., 2022)	tree	offline data	no assumption	1,600	15 states tree-MDP

Table 1. Our methods impose weaker assumptions on the data, and our empirical study covers more challenging tasks.

which is rather restrictive because many MDPs of interest then execute Monte Carlo methods inside the learned model. are not tree-structured. For example, in finite horizon MDPs Learning a high-fidelity model is, however, sometimes even considered in this paper, if two states at time  $t$  can transit to more challenging than evaluating the policy itself (Li, 2019). the same successor state at time  $t+1$ , then this MDP is not tree-structured. And the model prediction error can easily compound over tree-structured. Moreover, Mukherjee et al. (2022) require more steps during model rollouts (Wan et al., 2019) to directly approximate the transition function of the MDP. Nevertheless, if a good model could somehow be learned, our by counting, making it essentially a model-based approach. Our work still helps reduce the required rollouts when Monte Mukherjee et al. (2022), therefore, suffer from all canonical Monte Carlo is applied within the learned model.

challenges in model learning (Sutton, 1990; Sutton et al., 2008; Deisenroth & Rasmussen, 2011; Chua et al., 2018). Model-free of line evaluation. Model-free of line evaluation methods rely on learning other quantities for policy By contrast, we work on general MDPs without making any evaluation, including density ratio (a.k.a. marginalized im- assumption regarding their underlying structures, and we portance sampling ratio, Liu et al. (2018); Nachum et al. do not need to approximate the transition function. Our (2019); Li (2019); Xie et al. (2019); Zhang et al. (2020); approach is model-free. Zhong et al. (2022) adjust the be- Mousavi et al. (2020); Uehara et al. (2020); Yang et al. havior policy by encouraging under-sampled data. Their (2020)) and state-action value function (Harutyunyan et al., of line data, however, has to be complete trajectories gener- 2016; Munos et al., 2016; Farajtabar et al., 2018; Le et al., ated by known policies. In their experiments, they also 2019; Precup et al., 2000). But those learning processes require the policies for generating of line data to be similar bring in bias, either due to the misspecification of the func- to the target policy since they do not have any importance tion class or due to the complexity of optimization. Conse- sampling. By contrast, our method copes well with of line quently, the estimation they make is biased, and it is hard data in the form of incomplete segments from probably un- to quantify such bias without restrictive assumptions. known behavior policies that can be arbitrarily different More- our knowledge, the only practical way in general settings to over, there is no theoretical guarantee that the estimates made by Zhong et al. (2022) are certify that their estimation is indeed accurate is to compare unbiased or consistent. By contrast, our estimate is always those estimations with Monte Carlo estimations.

Other attempts for variance reduction in Monte Carlo evaluation mostly use control variates based on value functions (Zinkevich et al., 2006; White & Bowling, 2009; Jiang & Li, 2016). Such control variates can be integrated into our estimator, which we, however, save for future work. Notably, our work differs from the doubly robust method in Jiang & Li (2016) in that they assume the behavior policy is fixed and given while we use the fact that we have the freedom to choose a behavior policy in many settings. Moreover, to account for the stochasticity from the transition function, they require to learn a model of the MDP accurately, while we achieve this in a model-free way. Finally, they do not confirm a reduced variance compared with the on-policy estimator while we do. Efforts have been made to perform model selection with only of line data without explicitly learning a model as well (Paine et al., 2020; Kumar et al., 2021; Xie & Jiang, 2021; Zhang & Jiang, 2021). Those of line model selection methods, however, rarely have a correctness guarantee without restrictive assumptions. To summarize, if obtaining online data is entirely impossible, existing of line evaluation methods without using any online data might be the only choices. These include model-based methods and model-free methods augmented by of line model selection. However, in

Model-based of line evaluation. One straightforward way to exploit of line data for policy evaluation is to learn a model of the MDP first, probably with supervised learning (Jiang & Li, 2016; Paduraru, 2013; Zhang et al., 2021), and

many scenarios, it is practical to assume that a small amount of online data is available. If, in addition, evaluation correctness should be honored, then the improved Monte Carlo method in this work might be a better choice. Using offline data to help online model selection is previously explored by Konyushova et al. (2021). In particular, they use offline data to decide which policy, among a given set of policies, should be given priority to evaluate. When it comes to this then the estimation error divided by the average estimation error of the on-policy Monte Carlo estimator after the first episode. Thus, the normalized estimation error of the on-policy Monte Carlo estimator starts from

## 7. Empirical Results

In this section, we present empirical results comparing our methods against three baselines: (1) the canonical on-policy Monte Carlo estimator, (2) off-policy Monte Carlo estimator with behavior policy search (BPS, Hanna et al. (2017)), and (3) robust on-policy sampling (ROS, Zhong et al. (2022)). We do not implement ReVar (Mukherjee et al., 2022) because it will incur in infinite loops if the MDP is not tree-structured. Our method first learns a behavior policy with given offline data using Algorithm 1, then the PDIS Monte Carlo estimator (1) is used to estimate the performance of the target policy, where the learned behavior policy is used to interact with the environment. We call our method Offline Data Informed (ODI) algorithm. Our implementation is made publicly available to facilitate future research. Our method is superior in data requirements and applicability as summarized in Table 1.

**Gridworld:** We first conduct experiments with linear function approximation in Gridworld with  $m^3$  states, i.e., it is an  $m \times m$  grid with the time horizon also being  $m$ . Specifically, we use Gridworld with  $m^3 = 1;000$  and  $m^3 = 27;000$ . We use randomly generated reward functions and randomly generated target policies. The offline data is generated by selecting random actions on uniformly random state distribution. We report the normalized estimation error of the

<sup>1</sup><https://github.com/ShuzeLiu/Behavior-Policy-Design-for-Policy-Evaluation>

Figure 1. Results on Gridworld. The curves are averaged over 900 trials (30 target policies, each having 30 independent runs). The shaded regions denote standard errors and are invisible for some curves because they are too small.

As shown in Figure 1, our method outperforms baselines by a large margin. In particular, as shown by the dotted line, in Gridworld with size  $1;000$ , to achieve the same estimation error that the on-policy Monte Carlo estimator achieves with 250 steps, our methods only need around 50 steps. In Gridworld with size  $27;000$ , to achieve the same estimation error that on-policy Monte Carlo estimator achieves with 750 steps, our methods only need around 400 steps, saving more than 40% of online interactions. The improvement in environments with size  $= 27;000$  is smaller than environments with size  $= 1;000$  because the amount of offline data is the same for both environments, i.e., the offline data coverage is worse for the Gridworld with size  $= 27;000$ . In fact, the offline data coverage for the Gridworld with size  $= 1;000$  and size  $= 27;000$  are 62.5% and 2.3%, respectively. More experiment details are in Appendix B.1.

On-policy MC	Ours with 2.3% of line data coverage	Ours with 4.6% of line data coverage	Ours with 18.4% of line data coverage	BPG	ROS
300	150	90	60	300	300
600	330	180	120	540	540
1200	540	420	270	990	990

Table 2. The above table is an extension of Figure 1 by adding experiments with 18.4% data coverage for our algorithm in Gridworld with size  $= 27;000$ . Each number is the number of steps needed to achieve the same estimation accuracy that the naive Monte Carlo achieves with 300-600-1200 steps. All numbers are averaged from 900 different runs over a wide range of policies. Standard errors are visualized in Figure 1 of our paper and are invisible for some algorithm curves because they are too small.



Figure 2. Results on Mujoco environments. Each curve is averaged over 100 trials (30 target policies, each having 30 independent runs). The shaded regions denote standard errors and are invisible for some curves because they are too small.

	On-policy MC	Ours	BPG	ROS	Improvement in Saved Episodes
Ant	100	81	91	103	$(100-81)/(100-91)=211.1\%$
Hopper	100	54	89	100	$(100-54)/(100-89)=418.2\%$
I. Pendulum	100	72	103	99	$(100-72)/(100-99)=2800\%$
I. D. Pendulum	100	35	95	90	$(100-35)/(100-90)=650\%$
Walker	100	70	92	91	$(100-70)/(100-91)=333.3\%$

Table 3. Episodes needed to achieve the same of estimation accuracy that on-policy Monte Carlo achieves with 100 episodes.

We also show our algorithm scales with offline data. As we increase the data coverage in the Gridworld with size = 27;000 by adding more offline data generated from many different distributions, our method improves the saved samples from 55% =  $(1200 - 540)/1200$  to 77.5% =  $(1200 - 270)/1200$  in the last row of Table 2. By comparison, the best over all previous state-of-the-art algorithms only saved 7.5% =  $(1200 - 990)/1200$  samples and do not have a mechanism to use of offline data because they can only utilize online trajectory.

MuJoCo: We then conduct experiments with neural network function approximation in MuJoCo (Todorov et al., 2012) robot simulation tasks. Since our methods are designed for discrete action space, we discretize the MuJoCo action space. Details about action space discretization, target policy generation, and offline data generation are provided in Appendix B.2. We report the normalized estimator error in Figure 2, where our methods are consistently better than baselines. In particular, as shown by the dotted line in Figure 2 and Table 3, our methods need much fewer episodes (save up to 65% episodes) to achieve the estimation error that the on-policy Monte Carlo estimator achieves with 100 episodes. Recognizing episodes may have different lengths in MuJoCo, we also provide in Appendix B.2 a version of Figure 2 with the x-axis being steps, where our methods are still consistently better.

It is worth mentioning that all hyperparameters of our methods required to learn are tuned offline and are the same across all MuJoCo and Gridworld experiments.

## 8. Conclusion

Monte Carlo methods are the most dominant approach for evaluating a policy. The development and deployment of almost all RL algorithms, including of offline RL algorithms, implicitly or explicitly depend on Monte Carlo methods more or less. For example, when an RL researcher wants to plot a curve of the agent performance against training steps, Monte Carlo methods are usually the first choice. Our method improves the offline data efficiency of Monte Carlo evaluation while maintaining its unbiasedness by learning a tailored behavior policy from offline data. The two main contributions are the provably better closed-form behavior policy (Theorem 3) and its alternative representation (Theorem 4). Extending them to temporal difference learning (Sutton, 1988) is a possible future work.

## Acknowledgements

We thank Yuxin Chen for warm revising comments and Haifeng Xu and Zhengkun Xiao for insightful discussions. This work is supported in part by the US National Science Foundation under grants III-2128019 and SLES-2331904.

## Impact Statement

This paper advances the field of reinforcement learning and machine learning. There are many potential societal consequences of our work, none of which we feel must be especially highlighted here.

## References

- Antos, A., Szepesvári, C., and Munos, R. Learning near-optimal policies with bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning* 2008.
- Benjamin Melamed, R. Y. *Modern Simulation and Modeling* (Wiley Series in Probability and Statistics). Wiley-Interscience, 1998.
- Carpentier, A., Munos, R., and Antos, A. Adaptive strategy for stratified monte carlo sampling. *Journal of Machine Learning Research* 2015.
- Chen, J. and Jiang, N. Information-theoretic considerations in batch reinforcement learning. *Proceedings of International Conference on Machine Learning* 2019.
- Chervonyi, Y., Dutta, P., Trochim, P., Voicu, O., Paduraru, C., Qian, C., Karagozler, E., Davis, J. Q., Chippendale, R., Bajaj, G., et al. Semi-analytical industrial cooling system model for reinforcement learning. *arXiv preprint arXiv:2207.13131* 2022.
- Chua, K., Calandra, R., McAllister, R., and Levine, S. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. *Advances in Neural Information Processing Systems* 2018.
- Deisenroth, M. P. and Rasmussen, C. E. PILCO: A model-based and data-efficient approach to policy search. In *Proceedings of the International Conference on Machine Learning* 2011.
- Ernst, D., Geurts, P., and Wehenkel, L. Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research* 2005.
- Farajtabar, M., Chow, Y., and Ghavamzadeh, M. More robust doubly robust off-policy evaluation. *Proceedings of the International Conference on Machine Learning* 2018.
- Fu, J., Kumar, A., Nachum, O., Tucker, G., and Levine, S. D4RL: datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219* 2020.
- Fujimoto, S., Meger, D., and Precup, D. Off-policy deep reinforcement learning without exploration. *Proceedings of the International Conference on Machine Learning* 2019.
- Fujimoto, S., Meger, D., Precup, D., Nachum, O., and Guo, S. S. Why should i trust you, bellman? the bellman error is a poor replacement for value error. *arXiv preprint arXiv:2201.12417* 2022.
- Geweke, J. Antithetic acceleration of monte carlo integration in bayesian inference. *Journal of Econometrics* 1988.
- Gülçehre, Ç., Wang, Z., Novikov, A., Paine, T., Colmenarejo, S. G., Zolna, K., Agarwal, R., Merel, J., Mankowitz, D. J., Paduraru, C., Dulac-Arnold, G., Li, J., Norouzi, M., Hoffman, M., Heess, N., and de Freitas, N. RL unplugged: A collection of benchmarks for offline reinforcement learning. In *Advances in Neural Information Processing Systems* 2020.
- Hanna, J. P., Thomas, P. S., Stone, P., and Niekum, S. Data-efficient policy evaluation through behavior policy search. In *Proceedings of the International Conference on Machine Learning* 2017.
- Harutyunyan, A., Bellemare, M. G., Stepleton, T., and Munos, R. Q() with off-policy corrections. In *Proceedings of the International Conference on Algorithmic Learning Theory* 2016.
- Hesterberg, T. Weighted average importance sampling and defensive mixture distribution. *Technometrics* 1995.
- Huang, S., Dossa, R. F. J., Ye, C., Braga, J., Chakraborty, D., Mehta, K., and Araújo, J. G. Cleanrl: High-quality single-episode implementations of deep reinforcement learning algorithms. *Journal of Machine Learning Research* 2022.
- Jiang, N. and Li, L. Doubly robust off-policy value evaluation for reinforcement learning. *Proceedings of the International Conference on Machine Learning* 2016.
- Kahn, H. and Marshall, A. W. Methods of reducing sample size in monte carlo computations. *Journal of the Operations Research Society of America* 1953.
- Kakutani, S. Markoff process and the dirichlet problem. In *Proceedings of the Japan Academy* 1945.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations* 2015.
- Koller, D. and Friedman, N. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, 2009.
- Konyushova, K., Chen, Y., Paine, T., Gulcehre, C., Paduraru, C., Mankowitz, D. J., Denil, M., and de Freitas, N. Active offline policy selection. In *Advances in Neural Information Processing Systems* 2021.
- Kumar, A., Singh, A., Tian, S., Finn, C., and Levine, S. A work in progress for offline model-free robotic reinforcement learning. In *Proceedings of the Annual Conference on Robot Learning* 2021.

- Lange, S., Gabel, T., and Riedmiller, M. Batch reinforcement learning. *Reinforcement learning: State-of-the-art* 2012.
- Le, H. M., Voloshin, C., and Yue, Y. Batch policy learning under constraints. *Proceedings of the International Conference on Machine Learning* 2019.
- Levine, S., Kumar, A., Tucker, G., and Fu, J. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643* 2020.
- Li, L. A perspective on off-policy evaluation in reinforcement learning. *Frontiers of Computer Science* 2019.
- Liu, Q., Li, L., Tang, Z., and Zhou, D. Breaking the curse of horizon: In finite-horizon off-policy estimation. *Advances in Neural Information Processing Systems* 2018.
- Maei, H. R. Gradient temporal-difference learning algorithms. PhD thesis, University of Alberta, 2011.
- Mathieu, M., Ozair, S., Srinivasan, S., Gulcehre, C., Zhang, S., Jiang, R., Paine, T. L., Powell, Zolna, K., Schrittwieser, J., Choi, D., Georgiev, P., Toyama, D., Huang, A., Ring, R., Babuschkin, I., Ewalds, T., Bordbar, M., Henderson, S., Colmenarejo, S. G., van den Oord, A., Czarnecki, W. M., de Freitas, N., and Vinyals, O. AlphaStar unplugged: Large-scale offline reinforcement learning. *arXiv preprint arXiv:2308.03526* 2023.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M. A., Fidjeland, A., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. Human-level control through deep reinforcement learning. *Nature* 2015.
- Mousavi, A., Li, L., Liu, Q., and Zhou, D. Black-box off-policy estimation for finite-horizon reinforcement learning. *Proceedings of the International Conference on Learning Representations* 2020.
- Mukherjee, S., Hanna, J. P., and Nowak, R. D. Revar: Strengthening policy evaluation via reduced variance sampling. *Proceedings of the Conference in Uncertainty in Artificial Intelligence*, 2022.
- Munos, R. Error bounds for approximate policy iteration. *Proceedings of the International Conference on Machine Learning* 2003.
- Munos, R. and Szepesvári, C. Finite-time bounds for fitted value iteration. *Journal of Machine Learning Research* 2008.
- Munos, R., Stepleton, T., Harutyunyan, A., and Bellemare, M. Safe and efficient off-policy reinforcement learning. *In Advances in Neural Information Processing Systems* 2016.
- Nachum, O., Chow, Y., Dai, B., and Li, L. Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections. *Advances in Neural Information Processing Systems* 2019.
- O'Donoghue, B., Osband, I., Munos, R., and Mnih, V. The uncertainty bellman equation and exploration. *Proceedings of the International Conference on Machine Learning* 2018.
- Owen, A. B. Monte Carlo theory, methods and examples 2013.
- Paduraru, C. Off-policy evaluation in Markov decision processes. PhD thesis, McGill University, 2013.
- Paine, T. L., Paduraru, C., Michi, A., Gulcehre, C., Zolna, K., Novikov, A., Wang, Z., and de Freitas, N. Hyperparameter selection for offline reinforcement learning. *arXiv preprint arXiv:2007.09055* 2020.
- Precup, D., Sutton, R. S., and Singh, S. P. Eligibility traces for off-policy policy evaluation. *Proceedings of the International Conference on Machine Learning* 2000.
- Puterman, M. L. Markov decision processes: discrete stochastic dynamic programming. John Wiley & Sons, 2014.
- Rubinstein, R. Y. Simulation and the Monte Carlo Method. Wiley, 1981.
- Schrittwieser, J., Hubert, T., Mandhane, A., Barekatin, M., Antonoglou, I., and Silver, D. Online and offline reinforcement learning by planning with a learned model. *In Advances in Neural Information Processing Systems* 2021.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* 2017.
- Sherstan, C., Bennett, B., Young, K., Ashley, D. R., White, A., White, M., and Sutton, R. S. Directly estimating the variance of the  $Q$ -return using temporal-difference methods. *arXiv preprint arXiv:1801.08287* 2018.
- Sutton, R. S. Learning to predict by the methods of temporal differences. *Machine Learning* 1988.
- Sutton, R. S. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. *In Proceedings of the International Conference on Machine Learning* 1990.

- Sutton, R. S. and Barto, A. G. Reinforcement Learning: An Introduction (2nd Edition) MIT press, 2018.
- Sutton, R. S., McAllester, D. A., Singh, S. P., and Mansour, Y. Policy gradient methods for reinforcement learning with function approximation. In Advances in Neural Information Processing Systems 1999.
- Sutton, R. S., Szepesvári, C., Geramifard, A., and Bowling, M. H. Dyna-style planning with linear function approximation and prioritized sweeping. Proceedings of the Conference in Uncertainty in Artificial Intelligence 2008.
- Sutton, R. S., Mahmood, A. R., and White, M. An emphatic approach to the problem of off-policy temporal-difference learning. Journal of Machine Learning Research 2016.
- Tamar, A., Castro, D. D., and Mannor, S. Learning the variance of the reward-to-go. Journal of Machine Learning Research 2016.
- Thomas, P. S. Safe reinforcement learning PhD thesis, University of Massachusetts Amherst, 2015.
- Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. Proceedings of the International Conference on Intelligent Robots and Systems 2012.
- Uehara, M., Huang, J., and Jiang, N. Minimax weight and q-function learning for off-policy evaluation. Proceedings of the International Conference on Machine Learning 2020.
- Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P., Oh, J., Horgan, D., Kroiss, M., Danihelka, I., Huang, A., Sifre, L., Cai, T., Agapiou, J. P., Jaderberg, M., Vezhnevets, A. S., Leblond, R., Pohlen, T., Dalibard, V., Budden, D., Sulsky, Y., Molloy, J., Paine, T. L., Gehre, Ç., Wang, Z., Pfaff, T., Wu, Y., Ring, R., Yogatama, D., Wünsch, D., McKinney, K., Smith, O., Schaul, T., Lillicrap, T. P., Kavukcuoglu, K., Hassabis, D., Apps, C., and Silver, D. Grandmaster level in starcraft II using multi-agent reinforcement learning. Nature 2019.
- Wan, Y., Zaheer, M., White, A., White, M., and Sutton, R. S. Planning with expectation models. Proceedings of the International Joint Conference on Artificial Intelligence 2019.
- Watkins, C. J. and Dayan, P. Q-learning. Machine Learning 1992.
- Watkins, C. J. C. H. Learning from delayed rewards PhD thesis, King's College, Cambridge, 1989.
- White, M. and Bowling, M. H. Learning a value analysis tool for agent evaluation. Proceedings of the International Joint Conference on Artificial Intelligence 2009.
- Xie, T. and Jiang, N. Batch value-function approximation with only realizability. In Proceedings of the International Conference on Machine Learning 2021.
- Xie, T., Ma, Y., and Wang, Y. Towards optimal off-policy evaluation for reinforcement learning with marginalized importance sampling. Advances in Neural Information Processing Systems 2019.
- Yang, M., Nachum, O., Dai, B., Li, L., and Schuurmans, D. Off-policy evaluation via the regularized lagrangian. In Advances in Neural Information Processing Systems 2020.
- Zhang, M. R., Paine, T. L., Nachum, O., Paduraru, C., Tucker, G., Wang, Z., and Norouzi, M. Autoregressive dynamics models for offline policy evaluation and optimization. In Proceedings of the International Conference on Learning Representations 2021.
- Zhang, S. Breaking the deadly triad in reinforcement learning. PhD thesis, University of Oxford, 2022.
- Zhang, S. A new challenge in policy evaluation. Proceedings of the AAAI Conference on Artificial Intelligence 2023.
- Zhang, S. and Jiang, N. Towards hyperparameter-free policy selection for offline reinforcement learning. Advances in Neural Information Processing Systems 2021.
- Zhang, S., Liu, B., and Whiteson, S. GradientDICE: Rethinking generalized offline estimation of stationary values. In Proceedings of the International Conference on Machine Learning 2020.
- Zhong, R., Zhang, D., Schöler, L., Albrecht, S. V., and Hanna, J. P. Robust on-policy sampling for data-efficient policy evaluation in reinforcement learning. Advances in Neural Information Processing Systems 2022.
- Zinkevich, M., Bowling, M., Bard, N., Kan, M., and Billings, D. Optimal unbiased estimators for evaluating agent performance. In Proceedings of the AAAI Conference on Artificial Intelligence 2006.

## A. Proofs

### A.1. Proof of Lemma 1

Proof.

$$\begin{aligned}
 E_A [ (A)q(A) ] &= \sum_{a \in \mathcal{A}} \frac{p(a)}{p(a)} q(a) \\
 &= \sum_{a \in \mathcal{A}} q(a) \\
 &= \sum_{a \in \mathcal{A}} q(a) + \sum_{a \in \mathcal{A}} q(a) \cdot \mathbb{1}_{\{p(a)=0\}} \quad (2) \\
 &= \sum_{a \in \mathcal{A}} q(a) \\
 &= E_A [q(A)]:
 \end{aligned}$$

The intuition in the third equation is that the sample where  $p(a) = 0$ , i.e., this sample does not contribute to the expectation anyway.  $\square$

### A.2. Proof of Lemma 2

Proof.

For a given  $\alpha$  and  $q$ , define

$$A_+ \doteq \{a \in \mathcal{A} \mid q(a) \neq 0\}.$$

For any  $\alpha \in \mathcal{A}$ , we expand the variance as

$$\begin{aligned}
 V_A [ (A)q(A) ] &= E_A [ ((A)q(A))^2 ] - E_A^2 [ (A)q(A) ] \\
 &= E_A [ ((A)q(A))^2 ] - E_A^2 [q(A)] \quad (\text{Lemma 1}) \\
 &= \sum_{a \in \mathcal{A}} \frac{p(a)q^2(a)}{p(a)} E_A^2 [q(A)] \\
 &= \sum_{a \in \mathcal{A}_+} \frac{p(a)q^2(a)}{p(a)} E_A^2 [q(A)] \quad (p(a)q(a) = 0; \forall a \notin \mathcal{A}_+) \\
 &= \sum_{a \in \mathcal{A}_+} \frac{p(a)q^2(a)}{p(a)} E_A [q(A)]^2 \quad (2)
 \end{aligned}$$

The second term is a constant and is unrelated to  $\alpha$ . Solving the optimization problem (8) is, therefore, equivalent to solving

$$\min_{\alpha \in \mathcal{A}_+} \sum_{a \in \mathcal{A}_+} \frac{p(a)q^2(a)}{p(a)}. \quad (25)$$

Case 1:  $j_A + j = 0$

In this case, the variance is always 0, so any  $\alpha \in \mathcal{A}_+$  is optimal. In particular,  $\alpha = \frac{1}{A}$  is optimal.

Case 2:  $j_A + j > 0$

The definition of  $\alpha$  in (7) can be equivalently expressed, using contraposition, as

$$\alpha = \arg \min_{\alpha \in \mathcal{A}_+} \sum_{a \in \mathcal{A}_+} \frac{p(a)q^2(a)}{p(a)} \quad (a) > 0g:$$

The optimization problem (25) can then be equivalently written as

$$\min_{\alpha \in \mathcal{A}_+} \sum_{a \in \mathcal{A}_+} \frac{p(a)q^2(a)}{p(a)} \quad (26)$$

$$\text{s.t. } (a) > 0 \quad \forall a \in \mathcal{A}_+ :$$

If for some  $a_0 \in \mathcal{A}_+$  we have  $(a_0) < 1$ , then there must exist some  $a_1 \in \mathcal{A}_+$  such that  $(a_1) > 0$ . Since  $a_0$  does not contribute to the summation in the objective function (26), we can move the probability mass  $a_0$  to some other  $a_1 \in \mathcal{A}_+$  to increase  $(a_1)$  to further decrease the objective. In other words, any optimal solution (26) must put all its mass on  $\mathcal{A}_+$ . This motivates the following problem

$$\begin{aligned} \min_{z \in \mathcal{Z}(\mathcal{A}_+)} \quad & \sum_{a \in \mathcal{A}_+} \frac{z(a)q^2(a)}{z(a)} \\ \text{s.t. } \quad & z(a) > 0 \quad \forall a \in \mathcal{A}_+ : \end{aligned} \tag{27}$$

In particular, if  $z^*$  is an optimal solution to (27), then an optimal solution to (26) can be constructed as

$$(a) = \begin{cases} z^*(a) & a \in \mathcal{A}_+ \\ 0 & \text{otherwise} \end{cases} \tag{28}$$

Let  $R_{++} \doteq (0; +\infty)$ .

According to the Cauchy-Schwarz inequality, for any  $x \in R_{++}^{|\mathcal{A}_+|}$ , we have

$$\left( \sum_{a \in \mathcal{A}_+} \frac{z(a)q^2(a)}{z(a)} \right) \left( \sum_{a \in \mathcal{A}_+} \frac{(a)jq(a)j}{z(a)} \right) \geq \left( \sum_{a \in \mathcal{A}_+} (a)jq(a)j \right)^2$$

It can be easily verified that the equality holds for

$$z(a) \doteq \frac{(a)jq(a)j}{\sum_{b \in \mathcal{A}_+} (b)jq(b)j} > 0$$

Since  $\sum_{a \in \mathcal{A}_+} z(a) = 1$ , we conclude that  $z$  is an optimal solution to (27). An optimal solution to (8) can then be constructed according to (28). Making use of the fact that  $(a)jq(a)j = 0$  for  $a \notin \mathcal{A}_+$ , this can be equivalently expressed as

$$(a) = \frac{(a)jq(a)j}{\sum_{b \in \mathcal{A}_+} (b)jq(b)j};$$

which completes the proof. □

### A.3. Proof of Lemma 3

Proof. We start by showing  $\sum_{a \in \mathcal{A}_+} (a) = 1$ . Lemma 1 ensures that  $\sum_{a \in \mathcal{A}_+} (a) = 1$ . We now show that  $\sum_{a \in \mathcal{A}_+} (a) = 1$ . For any  $a \in \mathcal{A}_+$ , we have

$$\sum_{a \in \mathcal{A}_+} (a) \frac{(a)jq(a)j}{(a)jq(a)j} = \sum_{a \in \mathcal{A}_+} (a)jq(a)j$$

This indicates that

$$\sum_{a \in \mathcal{A}_+} (a)jq(a)j = 0$$

Since  $(a) \geq 0$  and all  $jq(a)j$  has the same sign, we must have

$$(a)jq(a)j = 0; \quad \forall a \in \mathcal{A}_+ : (a) = 0$$

This is exactly  $(a) = 0 \Rightarrow (a)jq(a)j = 0$ , yielding  $\sum_{a \in \mathcal{A}_+} (a) = 1$ . This completes the proof of  $\sum_{a \in \mathcal{A}_+} (a) = 1$ .

We now show the zero variance. When  $\pi(a) = 0$ , if  $q(a) = 0$ , we have  $\pi(a)q(a) = 0$ .

$$q(a) = \frac{(\pi(a))^j}{c}$$

and  $c > 0$  is a normalizing constant. Plugging to  $(\pi(a))q(a)$ , we get

$$(\pi(a))q(a) = \frac{(\pi(a))}{(\pi(a))} q(a) = \frac{(\pi(a))}{\frac{(\pi(a))^j}{c}} q(a) = c:$$

This means in this setting, with the optimal distribution, the random variable  $(\pi(a))q(a)$  is a constant function. Thus,

$$V_A((\pi(a))q(a)) = 0:$$

When  $\pi(a) = 0$ , if  $q(a) = 0$ , we have  $\pi(a)q(a) = 0$ .

$$q(a) = \frac{1}{|A|}$$

Plugging to  $(\pi(a))q(a)$ , we get

$$(\pi(a))q(a) = \frac{(\pi(a))}{(\pi(a))} q(a) = \frac{(\pi(a))q(a)}{1} = 0:$$

This shows  $(\pi(a))q(a)$  is also a constant. Thus,

$$V_A((\pi(a))q(a)) = 0:$$

The proof is similar for  $\pi(a) = 0$  and is thus omitted. □

#### A.4. Proof of Theorem 1

Proof. We proceed via induction. For  $t = T - 1$ , we have

$$\begin{aligned} E[G^{\text{PDIS}}(t:T-1) | S_t] &= E[R_{t+1} | S_t] = E[q_{t:t}(S_t; A_t) | S_t] \\ &= E_{A_t | (jS_t)} [q_{t:t}(S_t; A_t) | S_t] \quad (\text{Lemma 1}) \\ &= v_{t:t}(S_t); \end{aligned}$$

For  $t \in [T - 2]$ , we have

$$\begin{aligned} &E[G^{\text{PDIS}}(t:T-1) | S_t] \\ &= E[R_{t+1} + G^{\text{PDIS}}(t+1:T-1) | S_t] \\ &= E[R_{t+1} | S_t] + E[G^{\text{PDIS}}(t+1:T-1) | S_t] \\ &= E[R_{t+1} | S_t] + E_{A_t | (jS_t); S_{t+1}} [E[G^{\text{PDIS}}(t+1:T-1) | S_t; A_t; S_{t+1} | S_t] \quad (\text{Law of total expectation}) \\ &= E[R_{t+1} | S_t] + E_{A_t | (jS_t); S_{t+1}} [E[G^{\text{PDIS}}(t+1:T-1) | S_{t+1} | S_t] \quad (\text{Conditional independence and Markov property}) \\ &= E[R_{t+1} | S_t] + E_{A_t | (jS_t); S_{t+1}} [v_{t+1:t+1}(S_{t+1}) | S_t] \quad (\text{Inductive hypothesis}) \\ &= E_{A_t | (jS_t)} [q_{t:t}(S_t; A_t) | S_t] \quad (\text{Definition of } q_{t:t}) \\ &= E_{A_t | (jS_t)} [q_{t:t}(S_t; A_t) | S_t] \quad (\text{Lemma 1}) \\ &= v_{t:t}(S_t); \end{aligned}$$

which completes the proof. □

## A.5. Proof of Theorem 2

To prove Theorem 2, we rely on a recursive expression of the PDIS Monte Carlo estimator summarized by the following lemma.

Lemma 4 (Recursive Expression of Variance) For any  $\gamma \in [0, 1]$ , for  $t = T - 1$ ,

$$V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t:T-1}^1) \mid \mathbf{S}_t \right) = E_{A_t} \left[ \gamma^2 q_{\gamma,t}^2(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right] + v_{\gamma,t}^2(\mathbf{S}_t);$$

for  $t \in [T - 2]$ ,

$$\begin{aligned} & V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t:T-1}^1) \mid \mathbf{S}_t \right) \\ &= E_{A_t} \left[ \gamma^2 E_{S_{t+1}} \left[ V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid \mathbf{S}_t, \mathbf{A}_t \right) + \gamma r(\mathbf{S}_t; A_t) + q_{\gamma,t}^2(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right] \right. \\ & \quad \left. + v_{\gamma,t}^2(\mathbf{S}_t) \right]; \end{aligned}$$

Proof. When  $t = T - 1$ , we have

$$\begin{aligned} & V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t:T-1}^1) \mid \mathbf{S}_t \right) \tag{29} \\ &= E_{A_t} \left[ V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t:T-1}^1) \mid \mathbf{S}_t; A_t \right) \mid \mathbf{S}_t \right] + V_{A_t} \left( E \left[ G^{\text{PDIS}}(\gamma; \mathbf{s}_{t:T-1}^1) \mid \mathbf{S}_t; A_t \right] \mid \mathbf{S}_t \right) \tag{Law of total variance} \\ &= E_{A_t} \left[ \gamma^2 V \left( r(\mathbf{S}_t; A_t) + G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid \mathbf{S}_t; A_t \right) \mid \mathbf{S}_t \right. \\ & \quad \left. + V_{A_t} \left( E \left[ r(\mathbf{S}_t; A_t) + G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid \mathbf{S}_t; A_t \right] \mid \mathbf{S}_t \right) \right] \tag{Using (2)} \\ &= E_{A_t} \left[ \gamma^2 V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid \mathbf{S}_t; A_t \right) \mid \mathbf{S}_t \right] + V_{A_t} \left( E \left[ r(\mathbf{S}_t; A_t) + G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid \mathbf{S}_t; A_t \right] \mid \mathbf{S}_t \right) \\ & \quad \tag{Deterministic reward} \\ &= E_{A_t} \left[ \gamma^2 V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid \mathbf{S}_t; A_t \right) \mid \mathbf{S}_t \right] + V_{A_t} \left( \gamma q_{\gamma,t}(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right); \end{aligned}$$

Further decomposing the first term, we have

$$\begin{aligned} & V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid \mathbf{S}_t; A_t \right) \tag{30} \\ &= E_{S_{t+1}} \left[ V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid \mathbf{S}_t; A_t; S_{t+1} \right) \mid \mathbf{S}_t; A_t \right. \\ & \quad \left. + V_{S_{t+1}} \left( E \left[ G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid \mathbf{S}_t; A_t; S_{t+1} \right] \mid \mathbf{S}_t; A_t \right) \right] \tag{Law of total variance} \\ &= E_{S_{t+1}} \left[ V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid S_{t+1} \right) \mid \mathbf{S}_t; A_t \right] + V_{S_{t+1}} \left( E \left[ G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid S_{t+1} \right] \mid \mathbf{S}_t; A_t \right) \tag{Markov property} \\ &= E_{S_{t+1}} \left[ V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid S_{t+1} \right) \mid \mathbf{S}_t; A_t \right] + V_{S_{t+1}} \left( v_{\gamma,t+1}(S_{t+1}) \mid \mathbf{S}_t; A_t \right); \tag{Theorem 1} \end{aligned}$$

With  $v_{\gamma,t}$  defined in (11), plugging (30) back to (29) yields

$$\begin{aligned} & V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t:T-1}^1) \mid \mathbf{S}_t \right) \\ &= E_{A_t} \left[ \gamma^2 E_{S_{t+1}} \left[ V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid S_{t+1} \right) \mid \mathbf{S}_t; A_t \right] + v_{\gamma,t}(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right. \\ & \quad \left. + V_{A_t} \left( \gamma q_{\gamma,t}(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right) \right] \\ &= E_{A_t} \left[ \gamma^2 E_{S_{t+1}} \left[ V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid S_{t+1} \right) \mid \mathbf{S}_t; A_t \right] + v_{\gamma,t}(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right. \\ & \quad \left. + E_{A_t} \left[ \gamma^2 q_{\gamma,t}^2(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right] - \left( E_{A_t} \left[ \gamma q_{\gamma,t}(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right] \right)^2 \right] \\ &= E_{A_t} \left[ \gamma^2 E_{S_{t+1}} \left[ V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t+1:T-1}^1) \mid S_{t+1} \right) \mid \mathbf{S}_t; A_t \right] + v_{\gamma,t}(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right. \\ & \quad \left. + E_{A_t} \left[ \gamma^2 q_{\gamma,t}^2(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right] - v_{\gamma,t}^2(\mathbf{S}_t) \right]; \tag{Lemma 1} \end{aligned}$$

When  $t = T - 1$ , we have

$$\begin{aligned} V \left( G^{\text{PDIS}}(\gamma; \mathbf{s}_{t:T-1}^1) \mid \mathbf{S}_t \right) &= V \left( \gamma r(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right) \\ &= V \left( \gamma q_{\gamma,t}(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right) \\ &= E_{A_t} \left[ \gamma^2 q_{\gamma,t}^2(\mathbf{S}_t; A_t) \mid \mathbf{S}_t \right] - v_{\gamma,t}^2(\mathbf{S}_t); \end{aligned}$$

which completes the proof.  $\square$



We restate and present the main proof of Theorem 2.

Theorem 2 (Optimal Behavior Policy) For any  $\tau$  and  $s$ , the behavior policy  $\pi_\tau(s; a)$  defined above is an optimal solution to the following problem

$$\min_{\pi_\tau} V_{G^{\text{PDIS}}}(\pi_\tau; s) \mid S_\tau = s;$$

where  $\pi_\tau \doteq \pi_{\tau-1} \circ \pi_{\tau-2} \circ \dots \circ \pi_1$  and  $\pi_\tau(s; a) = 0$  if  $s \notin \mathcal{S}$  or  $a \notin \mathcal{A}(s)$ .

Proof. We proceed via induction. When  $\tau = 1$ , we have

$$\begin{aligned} & V_{G^{\text{PDIS}}}(\pi_1; s) \mid S_1 = s \\ &= V_{A_1}(\pi_1; s) \mid S_1 = s \\ &= V_{A_1}(\pi_1; s) \mid S_1 = s: \end{aligned}$$

The definition of  $\pi_1$  in (12) and Lemma 2 ensure that  $\pi_1$  is an optimal solution to

$$\min_{\pi_1} V_{G^{\text{PDIS}}}(\pi_1; s) \mid S_1 = s:$$

Now, suppose for some  $\tau \in [T-2]$ ,  $\pi_{\tau+1}$  is an optimal solution to

$$\min_{\pi_{\tau+1}} V_{G^{\text{PDIS}}}(\pi_{\tau+1}; s) \mid S_{\tau+1} = s:$$

To complete induction, we proceed to proving that  $\pi_\tau$  is an optimal solution to

$$\min_{\pi_\tau} V_{G^{\text{PDIS}}}(\pi_\tau; s) \mid S_\tau = s: \quad (31)$$

In the rest of this proof, we omit the domain  $\mathcal{S}_{\tau+1}$  for simplifying notations. For any  $s \in \mathcal{S}_\tau$ , we have

$$\begin{aligned} & V_{G^{\text{PDIS}}}(\pi_\tau; s) \mid S_\tau = s \\ &= E_{A_\tau} \left[ E_{S_{\tau+1}} \left[ V_{G^{\text{PDIS}}}(\pi_{\tau+1}; s) \mid S_{\tau+1} = s \mid S_\tau = s, A_\tau = a \right] + q_{\tau, \tau+1}^2(s; a) \right] \mid S_\tau = s \\ & \quad v_{\tau, \tau+1}^2(s) \quad \text{#} \quad \text{(By Lemma 4)} \\ & \stackrel{(a)}{=} E_{A_\tau} \left[ E_{S_{\tau+1}} \left[ \min_{\pi_{\tau+1}} V_{G^{\text{PDIS}}}(\pi_{\tau+1}; s) \mid S_{\tau+1} = s \mid S_\tau = s, A_\tau = a \right] + q_{\tau, \tau+1}^2(s; a) \right] \mid S_\tau = s \\ & \quad v_{\tau, \tau+1}^2(s) \quad \text{h} \quad \text{(Monotonically non-increasing in } \pi_{\tau+1} \text{)} \\ &= E_{A_\tau} \left[ E_{S_{\tau+1}} \left[ V_{G^{\text{PDIS}}}(\pi_{\tau+1}; s) \mid S_{\tau+1} = s \mid S_\tau = s, A_\tau = a \right] + q_{\tau, \tau+1}^2(s; a) \right] \mid S_\tau = s \\ & \quad v_{\tau, \tau+1}^2(s) \quad \text{(Inductive hypothesis)} \\ &= E_{A_\tau} \left[ E_{S_{\tau+1}} \left[ u_{\tau, \tau+1}(s; a) \mid S_\tau = s \right] + q_{\tau, \tau+1}^2(s; a) \right] \mid S_\tau = s \quad \text{(By (13))} \\ &= V_{A_\tau} \left[ E_{S_{\tau+1}} \left[ u_{\tau, \tau+1}(s; a) \mid S_\tau = s \right] + q_{\tau, \tau+1}^2(s; a) \right] \mid S_\tau = s \quad \text{(Definition of variance)} \\ &= V_{A_\tau} \left[ E_{S_{\tau+1}} \left[ u_{\tau, \tau+1}(s; a) \mid S_\tau = s \right] + q_{\tau, \tau+1}^2(s; a) \right] \mid S_\tau = s \quad \text{(Lemma 1 and } \tau \geq 1 \text{)} \\ & \stackrel{(b)}{=} E_{A_\tau} \left[ E_{S_{\tau+1}} \left[ u_{\tau, \tau+1}(s; a) \mid S_\tau = s \right]^2 + q_{\tau, \tau+1}^2(s; a) \right] \mid S_\tau = s \quad \text{(Non-negativity of variance)} \end{aligned}$$

According to the inductive hypothesis, the equality (a) can be achieved when  $\pi_{\tau+1} = \pi_{\tau+1}^*$ . According to the construction of  $\pi_\tau$  in (12) and Lemma 3, the equality (b) can be achieved when  $\pi_\tau = \pi_\tau^*$ . This suggests that  $\pi_\tau^*$  achieves the lower bound and is thus an optimal solution (31) which completes the induction and thus completes the proof.  $\square$

## A.6. Proof of Theorem 3

To prove the variance reduction property of the expression  $V_{G^{\text{PDIS}}(\hat{\pi}_{t:T-1}^1)} | S_t = s$ , the variance of the on-policy Monte Carlo estimator, in the form of a Bellman equation (Tamar et al., 2016; O'Donoghue et al., 2018; Sherstan et al., 2018). Define

$$r_{:t}(s; a) \doteq \left( r_{:t}(s; a) + \alpha_{:t}^2(s; a) - v_{:t}^2(s) \right) \mathbb{1}_{t \in [T-1]}; \quad (32)$$

$$q_{:t}(s; a) \doteq \begin{cases} r_{:t}(s; a) + \sum_{s^0; a^0} p(s^0; s; a) \mathbb{1}_{t+1}(a^0; s^0) q_{:t+1}(s^0; a^0) & \text{if } t \in [T-2] \\ r_{:t}(s; a) & \text{if } t = T-1 \end{cases}. \quad (33)$$

We have

Lemma 5 (Variance Equality)

$$V_{G^{\text{PDIS}}(\hat{\pi}_{t:T-1}^1)} | S_t = s = \sum_a \mathbb{1}_{t(a; s)} q_{:t}(s; a) - \mathbb{1}_{t; s}.$$

Proof. We proceed via induction. When  $t = T-1$ , we have

$$\begin{aligned} & V_{G^{\text{PDIS}}(\hat{\pi}_{t:T-1}^1)} | S_t \\ &= V_{A_t} (r_{:t}(S_t; A_t) | S_t) \\ &= V_{A_t} (r(S_t; A_t) | S_t) && \text{(By on-policy)} \\ &= V_{A_t} (q_{:t}(S_t; A_t) | S_t) \\ &= E_{A_t} \left[ \alpha_{:t}^2(S_t; A_t) | S_t \right] - v_{:t}^2(S_t) \\ &= \sum_a \mathbb{1}_{t(a; S_t)} q_{:t}(S_t; a); && \text{(By (33) and } \mathbb{1}_{T-1}(s; a) = 0 \text{)} \end{aligned}$$

For  $t \in [T-2]$ , we have

$$\begin{aligned} & V_{G^{\text{PDIS}}(\hat{\pi}_{t:T-1}^1)} | S_t \\ &= E_{A_t} \left[ E_{S_{t+1}} \left[ V_{G^{\text{PDIS}}(\hat{\pi}_{t+1:T-1}^1)} | S_{t+1} | S_t; A_t \right] + \alpha_{:t}^2(S_t; A_t) + r_{:t}(S_t; A_t) \right] - v_{:t}^2(S_t) \\ & && \text{(Lemma 4 and on-policy)} \\ &= \sum_a \mathbb{1}_{t(a; S_t)} \sum_{s^0} p(s^0; S_t; a) V_{G^{\text{PDIS}}(\hat{\pi}_{t+1:T-1}^1)} | S_{t+1} = s^0 + \mathbb{1}_{t; S_t; a} \\ &= \sum_a \mathbb{1}_{t(a; S_t)} \sum_{s^0} p(s^0; S_t; a) \sum_{a^0} \mathbb{1}_{t+1}(a^0; s^0) q_{:t+1}(s^0; a^0) + \mathbb{1}_{t; S_t; a} && \text{(Inductive hypothesis)} \\ &= \sum_a \mathbb{1}_{t(a; S_t)} q_{:t}(S_t; a); && \text{(By (33))} \end{aligned}$$

which completes the proof.  $\square$

Here, this is exactly the state-action value function of the target policy on the MDP w.r.t. to a new reward function. Manipulating (15) then yields

$$\begin{aligned} \hat{q}_{:t}(s; a) &= \sum_{s^0} p(s^0; s; a) \sum_{a^0} \mathbb{1}_{t+1}(a^0; s^0) q_{:t+1}(s^0; a^0) + \mathbb{1}_{t; s; a} + \alpha_{:t}^2(s; a) \\ &= q_{:t}(s; a) + v_{:t}^2(s); \end{aligned} \quad (34)$$

Now, we restate and present the main proof of Theorem 3.

Theorem 3 (Variance Reduction) For any  $\alpha$  and,

$$V_{G^{\text{PDIS}}(\hat{\pi}_{t:T-1}^1)} | S_t = s$$

$$V G^{\text{PDIS}}(\hat{\pi}_{t:T-1}^1) | S_t = s_t(s):$$

To denote  $v_t(s)$ , first denote  $c_t(s) =$

$$P_{a_t} (aj_s) \hat{q}_{t,t}(s; a) - P_{a_t} (aj_s) \overline{\hat{q}_{t,t}(s; a)}^2:$$

Then we denote  $\dot{c}_t(s) \doteq c_t(s)$  for  $t = T-1$  and otherwise

$$\dot{c}_t(s) \doteq c_t(s) + E_{A_t} \hat{\pi}_t^2 E_{S_{t+1}} [\hat{\pi}_{t+1}(S_{t+1}) | s; A_t] : \quad (19)$$

Proof. We proceed via induction. For  $t = T-1$ , we have

$$\begin{aligned} & V G^{\text{PDIS}}(\hat{\pi}_{t:T-1}^1) | S_t \\ &= E_{A_t} \hat{\pi}_t^2 \hat{q}_{t,t}^2(S_t; A_t) | S_t - v_{t,t}^2(S_t) \quad (\text{Lemma 4}) \\ &= E_{A_t} \hat{\pi}_t^2 \hat{q}_{t,t}(S_t; A_t) | S_t - v_{t,t}^2(S_t) \quad (\text{Definition of } \hat{q} \text{ (14)}) \\ &= V_{A_t} \hat{\pi}_t^2 \hat{q}_{t,t}(S_t; A_t) | S_t + E_{A_t}^2 \hat{\pi}_t^2 \hat{q}_{t,t}(S_t; A_t) | S_t - v_{t,t}^2(S_t) \\ & \quad (\text{Definition of variance and non-negativity of } \hat{q}) \\ &= V_{A_t} \hat{\pi}_t^2 \hat{q}_{t,t}(S_t; A_t) | S_t + \sum_a \hat{\pi}_t(a | S_t) \hat{q}_{t,t}(S_t; a)^2 - v_{t,t}^2(S_t) \quad (\text{Lemma 1}) \\ &= \sum_a \hat{\pi}_t(a | S_t) \hat{q}_{t,t}(S_t; a)^2 - v_{t,t}^2(S_t) \quad (\text{Definition of } \hat{\pi} \text{ (18) and Lemma 3}) \\ &= \sum_a \hat{\pi}_t(a | S_t) \hat{q}_{t,t}(S_t; a) + \sum_a \hat{\pi}_t(a | S_t) \hat{q}_{t,t}(S_t; a)^2 - \sum_a \hat{\pi}_t(a | S_t) \hat{q}_{t,t}(S_t; a) - v_{t,t}^2(S_t) \\ &= V G^{\text{PDIS}}(\hat{\pi}_{t:T-1}^1) | S_t + \sum_a \hat{\pi}_t(a | S_t) \hat{q}_{t,t}(S_t; a)^2 - \sum_a \hat{\pi}_t(a | S_t) \hat{q}_{t,t}(S_t; a) \quad (\text{By (34) and Lemma 5}) \\ &= V G^{\text{PDIS}}(\hat{\pi}_{t:T-1}^1) | S_t - \dot{c}_t(S_t): \quad (\text{Definition of } \dot{c} \text{ (19)}) \end{aligned}$$

For  $t \in [T-2]$ , we have

$$\begin{aligned} & V G^{\text{PDIS}}(\hat{\pi}_{t:T-1}^1) | S_t \\ &= E_{A_t} \hat{\pi}_t^2 E_{S_{t+1}} V G^{\text{PDIS}}(\hat{\pi}_{t+1:T-1}^1) | S_{t+1} | S_t; A_t + \hat{\pi}_t(S_t; A_t) + \hat{q}_{t,t}^2(S_t; A_t) | S_t \\ & \quad - v_{t,t}^2(S_t) \quad (\text{Lemma 4}) \\ &= E_{A_t} \hat{\pi}_t^2 E_{S_{t+1}} \sum_a \hat{\pi}_{t+1}(a | S_{t+1}) \hat{q}_{t,t+1}(S_{t+1}; a) | S_t; A_t + \hat{\pi}_t(S_t; A_t) \\ & \quad + \hat{q}_{t,t}^2(S_t; A_t) | S_t - v_{t,t}^2(S_t) - E_{A_t} \hat{\pi}_t^2 E_{S_{t+1}} [\hat{\pi}_{t+1}(S_{t+1}) | S_t; A_t] \quad (\text{Inductive hypothesis and Lemma 5}) \\ &= E_{A_t} \hat{\pi}_t^2 \hat{q}_{t,t}(S_t; A_t) + v_{t,t}^2(S_t) | S_t - v_{t,t}^2(S_t) - E_{A_t} \hat{\pi}_t^2 E_{S_{t+1}} [\hat{\pi}_{t+1}(S_{t+1}) | S_t; A_t] \\ & \quad (\text{Definition of } \hat{q} \text{ (33)}) \\ &= E_{A_t} \hat{\pi}_t^2 \hat{q}_{t,t}(S_t; A_t) | S_t - v_{t,t}^2(S_t) - E_{A_t} \hat{\pi}_t^2 E_{S_{t+1}} [\hat{\pi}_{t+1}(S_{t+1}) | S_t; A_t] \quad (\text{Definition of } \hat{q} \text{ (34)}) \\ &= V_{A_t} \hat{\pi}_t^2 \hat{q}_{t,t}(S_t; A_t) | S_t + E_{A_t}^2 \hat{\pi}_t^2 \hat{q}_{t,t}(S_t; A_t) | S_t - v_{t,t}^2(S_t) \\ & \quad - E_{A_t} \hat{\pi}_t^2 E_{S_{t+1}} [\hat{\pi}_{t+1}(S_{t+1}) | S_t; A_t] \quad (\text{Definition of variance and non-negativity of } \hat{q}) \\ &= V_{A_t} \hat{\pi}_t^2 \hat{q}_{t,t}(S_t; A_t) | S_t + \sum_a \hat{\pi}_t(a | S_t) \hat{q}_{t,t}(S_t; a)^2 - v_{t,t}^2(S_t) \end{aligned}$$

$$\begin{aligned}
 & E_{A_t, \Lambda_t} \left[ \sum_{a \in \mathcal{A}} q_{t+1}(S_{t+1}; a) \left( \sum_{a' \in \mathcal{A}} q_{t+1}(S_{t+1}; a') v_{t+1}^2(S_{t+1}; a') - v_{t+1}^2(S_{t+1}; a) \right) \mid S_t; A_t \right] && \text{(Lemma 1)} \\
 = & \sum_{a \in \mathcal{A}} q_{t+1}(S_{t+1}; a) \left( \sum_{a' \in \mathcal{A}} q_{t+1}(S_{t+1}; a') v_{t+1}^2(S_{t+1}; a') - v_{t+1}^2(S_{t+1}; a) \right) E_{A_t, \Lambda_t} \left[ \sum_{a' \in \mathcal{A}} q_{t+1}(S_{t+1}; a') \mid S_t; A_t \right] && \text{(Definition of } \Lambda \text{ (18) and Lemma 3)} \\
 = & \sum_{a \in \mathcal{A}} q_{t+1}(S_{t+1}; a) v_{t+1}^2(S_{t+1}; a) + \sum_{a \in \mathcal{A}} q_{t+1}(S_{t+1}; a) \left( \sum_{a' \in \mathcal{A}} q_{t+1}(S_{t+1}; a') v_{t+1}^2(S_{t+1}; a') - v_{t+1}^2(S_{t+1}; a) \right) && \\
 = & V \text{G}^{\text{PDIS}}(r_{t:T}, \mathbf{1}) \mid S_t + \sum_{a \in \mathcal{A}} q_{t+1}(S_{t+1}; a) \left( \sum_{a' \in \mathcal{A}} q_{t+1}(S_{t+1}; a') v_{t+1}^2(S_{t+1}; a') - v_{t+1}^2(S_{t+1}; a) \right) && \\
 = & V \text{G}^{\text{PDIS}}(r_{t:T}, \mathbf{1}) \mid S_t + \sum_{a \in \mathcal{A}} q_{t+1}(S_{t+1}; a) v_{t+1}^2(S_{t+1}; a) && \text{(By (34) and Lemma 5)} \\
 = & V \text{G}^{\text{PDIS}}(r_{t:T}, \mathbf{1}) \mid S_t + \sum_{a \in \mathcal{A}} q_{t+1}(S_{t+1}; a) v_{t+1}^2(S_{t+1}; a) && \text{(Definition of (19))}
 \end{aligned}$$

□

#### A.7. Proof of Theorem 4

Proof. For  $t = T - 1$ , we have

$$\begin{aligned}
 \hat{q}_{t+1}(s; a) &= q_{t+1}^2(s; a) && \text{(Definition of } \hat{q}_{t+1} \text{ (14))} \\
 &= \hat{r}_{t+1}(s; a) && \text{(By } q_{t+1}(s; a) = r(s; a) \text{ and Theorem 4)}
 \end{aligned}$$

For  $t \in [T - 2]$ , we have

$$\begin{aligned}
 & \hat{q}_{t+1}(s; a) && \\
 = & \hat{r}_{t+1}(s; a) + v_{t+1}^2(s; a) && \text{(By (34))} \\
 = & \hat{r}_{t+1}(s; a) + v_{t+1}^2(s; a) + \sum_{s^0, a^0} p(s^0; s; a) \left( \hat{q}_{t+1}(s^0; a^0) - v_{t+1}^2(s^0; a^0) \right) && \text{(Definition of } \hat{q} \text{ (33))} \\
 = & \hat{r}_{t+1}(s; a) + v_{t+1}^2(s; a) + \sum_{s^0, a^0} p(s^0; s; a) \left( \hat{q}_{t+1}(s^0; a^0) - v_{t+1}^2(s^0; a^0) \right) && \\
 = & \hat{r}_{t+1}(s; a) + v_{t+1}^2(s; a) + \sum_{s^0, a^0} p(s^0; s; a) \left( \hat{q}_{t+1}(s^0; a^0) - v_{t+1}^2(s^0; a^0) \right) && \text{(By (34))} \\
 = & \hat{r}_{t+1}(s; a) + q_{t+1}^2(s; a) + \sum_{s^0, a^0} p(s^0; s; a) v_{t+1}^2(s^0; a^0) + \sum_{s^0, a^0} p(s^0; s; a) \left( \hat{q}_{t+1}(s^0; a^0) - v_{t+1}^2(s^0; a^0) \right) && \text{(Definition of } \hat{r} \text{ (32))} \\
 = & (E[v_{t+1}(S_{t+1}) \mid S_t = s; A_t = a])^2 + q_{t+1}^2(s; a) + \sum_{s^0, a^0} p(s^0; s; a) \left( \hat{q}_{t+1}(s^0; a^0) - v_{t+1}^2(s^0; a^0) \right) && \text{(Definition of (11))} \\
 = & (q_{t+1}(s; a) - r(s; a))^2 + q_{t+1}^2(s; a) + \sum_{s^0, a^0} p(s^0; s; a) \left( \hat{q}_{t+1}(s^0; a^0) - v_{t+1}^2(s^0; a^0) \right) && \\
 = & 2r(s; a)q_{t+1}(s; a) - r^2(s; a) + \sum_{s^0, a^0} p(s^0; s; a) \left( \hat{q}_{t+1}(s^0; a^0) - v_{t+1}^2(s^0; a^0) \right) && \\
 = & \hat{r}_{t+1}(s; a) + \sum_{s^0, a^0} p(s^0; s; a) \left( \hat{q}_{t+1}(s^0; a^0) - v_{t+1}^2(s^0; a^0) \right); && \text{(By Theorem 4)}
 \end{aligned}$$

which completes the proof. □

#### A.8. Proof of Theorem 5

Proof. We first derive an important equality

$$E_{A_t, \Lambda_t} \left[ \sum_{a \in \mathcal{A}} q_{t+1}^2(S_{t+1}; a) \mid S_t; A_t \right]$$

$$\begin{aligned}
 &= \sum_a \frac{q_t^2(a|S_t)}{\hat{\lambda}_t^+(a|S_t)} \hat{q}_{t,t}(S_t; a) \\
 &= \sum_a \frac{q_t^2(a|S_t)}{P_{b \sim \pi_t(b|S_t)} \frac{q_{t,t}^+(S_t; a)}{q_{t,t}^+(S_t; b)}} \hat{q}_{t,t}(S_t; a) \quad (\text{by (23)}) \\
 &= \sum_a \frac{q_t(a|S_t)}{q_{t,t}^+(S_t; a)} \sum_a \frac{q_{t,t}^+(S_t; a)}{q_{t,t}^+(S_t; a)} \hat{q}_{t,t}(S_t; a) \\
 &= \sum_a \frac{q_t(a|S_t)}{q_{t,t}(S_t; a)} \sum_a \frac{q_{t,t}(S_t; a)}{q_{t,t}^+(S_t; a)} \hat{q}_{t,t}(S_t; a) \quad (\text{By (22)}) \quad (35)
 \end{aligned}$$

We proceed via induction. For  $t = T - 1$ , we have

$$\begin{aligned}
 &V G^{\text{PDIS}}(\hat{\lambda}_{t:T}^{+1}) \mid S_t \\
 &= E_{A_t} \sum_a \frac{q_t^2(a|S_t)}{h} \hat{q}_{t,t}(S_t; A_t) \mid S_t \quad v_{t,t}^2(S_t) \quad (\text{Lemma 4}) \\
 &= E_{A_t} \sum_a \frac{q_t^2(a|S_t)}{h} \hat{q}_{t,t}(S_t; A_t) \mid S_t \quad v_{t,t}^2(S_t) \quad (\text{Definition of } \hat{q} \text{ (14)}) \\
 &= \sum_a \frac{q_t(a|S_t)}{q_{t,t}(S_t; a)} \sum_a \frac{q_{t,t}(S_t; a)}{q_{t,t}^+(S_t; a)} \sum_a \frac{q_{t,t}^+(S_t; a)}{q_{t,t}^+(S_t; a)} \hat{q}_{t,t}(S_t; a) \quad v_{t,t}^2(S_t) \quad (\text{By (35)}) \\
 &= \sum_a \frac{q_t(a|S_t)}{q_{t,t}(S_t; a)} \hat{q}_{t,t}(S_t; a) + \sum_a \frac{q_t(a|S_t)}{q_{t,t}(S_t; a)} \sum_a \frac{q_{t,t}(S_t; a)}{q_{t,t}^+(S_t; a)} \sum_a \frac{q_{t,t}^+(S_t; a)}{q_{t,t}^+(S_t; a)} \hat{q}_{t,t}(S_t; a) \\
 &= \sum_a \frac{q_t(a|S_t)}{q_{t,t}(S_t; a)} \hat{q}_{t,t}(S_t; a) \quad v_{t,t}^2(S_t) \\
 &= V G^{\text{PDIS}}(\hat{\lambda}_{t:T}^{+1}) \mid S_t \quad \sum_a \frac{q_t(a|S_t)}{q_{t,t}(S_t; a)} \sum_a \frac{q_{t,t}(S_t; a)}{q_{t,t}^+(S_t; a)} \sum_a \frac{q_{t,t}^+(S_t; a)}{q_{t,t}^+(S_t; a)} \hat{q}_{t,t}(S_t; a) \\
 &= V G^{\text{PDIS}}(\hat{\lambda}_{t:T}^{+1}) \mid S_t \quad v_{t,t}^2(S_t) \quad (\text{By (34) and Lemma 5}) \\
 &= V G^{\text{PDIS}}(\hat{\lambda}_{t:T}^{+1}) \mid S_t \quad v_{t,t}^2(S_t) \quad (\text{Definition of } v^+ \text{ (24)})
 \end{aligned}$$

For  $t \in [T - 2]$ , we have

$$\begin{aligned}
 &V G^{\text{PDIS}}(\hat{\lambda}_{t:T}^{+1}) \mid S_t \\
 &= E_{A_t} \sum_a \frac{q_t^2(a|S_t)}{h} E_{S_{t+1}} V G^{\text{PDIS}}(\hat{\lambda}_{t+1:T}^{+1}) \mid S_{t+1} \mid S_t; A_t + \sum_a \frac{q_t^2(a|S_t)}{h} \hat{q}_{t,t}(S_t; A_t) \mid S_t \\
 &\quad v_{t,t}^2(S_t) \quad (\text{Lemma 4}) \\
 &= E_{A_t} \sum_a \frac{q_t^2(a|S_t)}{h} E_{S_{t+1}} \sum_a \frac{q_{t+1}(a|S_{t+1})}{h} \hat{q}_{t+1,t+1}(S_{t+1}; a) \mid S_t; A_t + \sum_a \frac{q_t^2(a|S_t)}{h} \hat{q}_{t,t}(S_t; A_t) \\
 &\quad + \sum_a \frac{q_t^2(a|S_t)}{h} \hat{q}_{t,t}(S_t; A_t) \mid S_t \quad v_{t,t}^2(S_t) \quad E_{A_t} \sum_a \frac{q_t^2(a|S_t)}{h} E_{S_{t+1}} \sum_a \frac{q_{t+1}(a|S_{t+1})}{h} \hat{q}_{t+1,t+1}(S_{t+1}) \mid S_t; A_t \quad (\text{Inductive hypothesis and Lemma 5}) \\
 &= E_{A_t} \sum_a \frac{q_t^2(a|S_t)}{h} \hat{q}_{t,t}(S_t; A_t) + v_{t,t}^2(S_t) \mid S_t \quad v_{t,t}^2(S_t) \quad E_{A_t} \sum_a \frac{q_t^2(a|S_t)}{h} E_{S_{t+1}} \sum_a \frac{q_{t+1}(a|S_{t+1})}{h} \hat{q}_{t+1,t+1}(S_{t+1}) \mid S_t; A_t \\
 &\quad (\text{Definition of } \hat{q} \text{ (33)}) \\
 &= E_{A_t} \sum_a \frac{q_t^2(a|S_t)}{h} \hat{q}_{t,t}(S_t; A_t) \mid S_t \quad v_{t,t}^2(S_t) \quad E_{A_t} \sum_a \frac{q_t^2(a|S_t)}{h} E_{S_{t+1}} \sum_a \frac{q_{t+1}(a|S_{t+1})}{h} \hat{q}_{t+1,t+1}(S_{t+1}) \mid S_t; A_t \\
 &\quad (\text{Definition of } \hat{q} \text{ (15)}) \\
 &= \sum_a \frac{q_t(a|S_t)}{q_{t,t}(S_t; a)} \sum_a \frac{q_{t,t}(S_t; a)}{q_{t,t}^+(S_t; a)} \sum_a \frac{q_{t,t}^+(S_t; a)}{q_{t,t}^+(S_t; a)} \hat{q}_{t,t}(S_t; a) \quad v_{t,t}^2(S_t)
 \end{aligned}$$

$$= \sum_{\mathbf{a}} \mathbb{E}_{A_t \sim \pi_t^+} \left[ \mathbb{E}_{S_{t+1} \sim p_{t+1}^+} (j(S_{t+1}) | S_t; A_t) \right] \quad (\text{By (35)})$$

$$= \sum_{\mathbf{a}} \mathbb{E}_{S_t} \left[ \frac{1}{\pi_t(\mathbf{a} | S_t)} \mathbb{E}_{S_{t+1} \sim p_{t+1}^+} (j(S_{t+1}) | S_t; \mathbf{a}) \right] + \sum_{\mathbf{a}} \mathbb{E}_{S_t} \left[ \frac{1}{\pi_t(\mathbf{a} | S_t)} \mathbb{E}_{S_{t+1} \sim p_{t+1}^+} (j(S_{t+1}) | S_t; \mathbf{a}) \right]$$

$$= V + \sum_{\mathbf{a}} \mathbb{E}_{S_t} \left[ \frac{1}{\pi_t(\mathbf{a} | S_t)} \mathbb{E}_{S_{t+1} \sim p_{t+1}^+} (j(S_{t+1}) | S_t; \mathbf{a}) \right] - \sum_{\mathbf{a}} \mathbb{E}_{S_t} \left[ \frac{1}{\pi_t(\mathbf{a} | S_t)} \mathbb{E}_{S_{t+1} \sim p_{t+1}^+} (j(S_{t+1}) | S_t; \mathbf{a}) \right] \quad (\text{By (34) and Lemma 5})$$

$$= V + \sum_{\mathbf{a}} \mathbb{E}_{S_t} \left[ \frac{1}{\pi_t(\mathbf{a} | S_t)} \mathbb{E}_{S_{t+1} \sim p_{t+1}^+} (j(S_{t+1}) | S_t; \mathbf{a}) \right] \quad (\text{Definition of } \pi_t^+ \text{ (24)})$$

□

## B. Experiment Details

### B.1. GridWorld

For a Gridworld with size  $n$ , its width, height, and time horizon  $T$  are all set to  $n$ . There are four possible actions: up, down, left, and right. After taking an action, the agent has a 0.9 probability of moving accordingly and a 0.1 probability of moving uniformly at random. If the agent runs into a boundary, the agent stays in its current location. The reward function  $r(s; a)$  is randomly generated and fixed after generation. We normalize the rewards across all  $(s; a)$  such that  $\max_{s; a} r(s; a) = 1$ . We consider a set of randomly generated target policies. The ground truth policy performance is estimated using the on-policy Monte Carlo method by running each target policy for  $10^6$  episodes. We test two different sizes of the Gridworld with a number of 1;000 and 27;000 states. The offline dataset contains  $m = 10^5$  randomly generated tuples. For a Gridworld of size  $n$ , the total amount of possible  $(s; t; a; r; s')$  tuples is  $n \cdot n \cdot n \cdot 4 \cdot 4 = 16n^3$ . The offline data coverages for the Gridworld of size 1;000 and 27;000 are then 62.5% and 2.3%.

We use a one-hot vector representing the position of the agent and a real number representing the current time step as features for the state. We execute Algorithm 1 to approximate function  $r$ ,  $q$ , and  $\hat{q}$ . As shown in Algorithm 1, we train  $r$  using supervised learning by batch stochastic gradient descent. We train  $q$  and  $\hat{q}$  using fitted  $Q$ -learning. We split the offline data into a training set and a test set. We tune all hyperparameters offline based on the supervised learning loss and fitted  $Q$ -learning loss on the test set. With the Adam optimizer (Kingma & Ba, 2015), we search the learning rates in  $2^{-20}; 2^{-18}; \dots; 2^0$  to minimize the loss on the offline data and use the learning rate  $2^{-10}$  on all learning processes. For the behavior policy search (BPS, Hanna et al. (2017)) and robust on-policy sampling (ROS, Zhong et al. (2022)) algorithms, we use the reported parameters from Hanna et al. (2017) and Zhong et al. (2022), since it is not clear how to do hyperparameter tuning for BPS and ROS with only offline data.

### B.2. MuJoCo

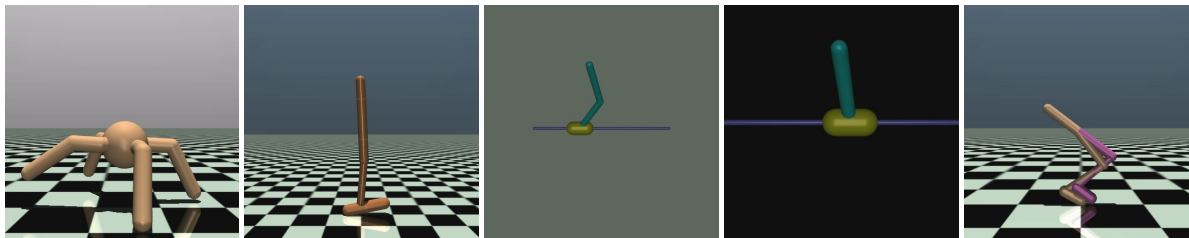


Figure 3. MuJoCo (Todorov et al., 2012) robot simulation tasks. MuJoCo is a physics engine for robotics simulation and contains various stochastic environments. The goal in each environment is to control a robot to achieve different behaviors such as walking, jumping, and balancing. Environments from the left to the right are Ant, Hopper, InvertedDoublePendulum, InvertedPendulum, and Walker. We conducted experiments on those five environments with results reported in Section 7.

Figure 3 is an introduction to the MuJoCo environments. We construct 150 policies (30 policies in each environment) with a wide range of performance using the proximal policy optimization (PPO) algorithm (Schulman et al., 2017) and the default PPO implementation in Huang et al. (2022). Since our methods are designed for discrete action space, we discretize the first dimension of MuJoCo action space in our experiments. The remaining dimensions are controlled by the PPO policy and are deemed as part of the environment. We run each compared algorithm 30 times for each policy and compute the average and standard error to plot curves in Figure 2. To generate offline data, we add different levels of noise to the target policy and run noisy target policies for 2000 episodes. The noise is in the form of a uniformly random policy, and its weight is uniformly randomly sampled from  $(0; 0.1]$ . This data generation process simulates the data generated during the training of a policy. Notably, compared with previous works, we do not need data to be complete trajectories or generated by known policies. We leave the investigation of entirely irrelevant offline data in the MuJoCo domain for future work. Our algorithm is robust on hyperparameters. All learning rates in Algorithm 1 are tuned offline and are the same  $2^{-10}$  across all MuJoCo and Gridworld experiments.

In MuJoCo, the episode length varies because of stochasticity in policies and environments. Because the length of each episode is not fixed, episodes in off-policy estimation may be longer than episodes in on-policy estimation. In the main text, we use episodes instead of steps as the  $x$ -axis mainly to improve readability. Because after running 100 steps, we might already have a good estimate for a target policy with a length of 10 but may still not finish a single episode for a target

