
Multi-step Greedy Reinforcement Learning Algorithms

Manan Tomar^{*1} Yonathan Efroni^{*2} Mohammad Ghavamzadeh³

Abstract

Multi-step greedy policies have been extensively used in model-based reinforcement learning (RL), both when a model of the environment is available (e.g., in the game of Go) and when it is learned. In this paper, we explore their benefits in model-free RL, when employed using multi-step dynamic programming algorithms: κ -Policy Iteration (κ -PI) and κ -Value Iteration (κ -VI). These methods iteratively compute the next policy (κ -PI) and value function (κ -VI) by solving a surrogate decision problem with a shaped reward and a smaller discount factor. We derive model-free RL algorithms based on κ -PI and κ -VI in which the surrogate problem can be solved by any discrete or continuous action RL method, such as DQN and TRPO. We identify the importance of a hyperparameter that controls the extent to which the surrogate problem is solved and suggest a way to set this parameter. When evaluated on a range of Atari and MuJoCo benchmark tasks, our results indicate that for the right range of κ , our algorithms outperform DQN and TRPO. This shows that our multi-step greedy algorithms are general enough to be applied over any existing RL algorithm and can significantly improve its performance.

1. Introduction

Reinforcement learning (RL) algorithms solve sequential decision-making problems through repeated interaction with the environment. By incorporating deep neural networks into RL algorithms, the field has recently witnessed remarkable empirical success (e.g., Mnih et al., 2015; Lillicrap et al., 2015; Levine et al., 2016; Silver et al., 2017). Much of this success has been achieved by model-free RL algorithms, such as Q-learning and policy gradients. These algorithms are known to suffer from high variance in their

estimations (Greensmith et al., 2004) and to have difficulties in handling function approximation (e.g., Thrun & Schwartz, 1993; Baird, 1995; Van Hasselt et al., 2016; Lu et al., 2018). These issues are intensified in decision problems with long horizon, i.e., when the discount factor, γ , is large. Although using smaller values of γ addresses the discount factor-dependent issues and leads to more stable algorithms (Patrik & Scherrer, 2009; Jiang et al., 2015), it does not come for free, as the algorithm may return a *biased* solution, i.e., it may not converge to an optimal (or good) solution for the original decision problem (the one with larger value of γ).

Efroni et al. (2018a) recently proposed another approach to mitigate the γ -dependant instabilities in RL in which they study multi-step greedy versions of the well-known Dynamic Programming (DP) algorithms: Policy Iteration (PI) and Value Iteration (VI) (Bertsekas & Tsitsiklis, 1996). They also proposed an alternative formulation of the multi-step greedy policy, called κ -greedy policy, and studied the convergence of the resulted PI and VI algorithms: κ -PI and κ -VI. These algorithms iteratively solve a $\gamma\kappa$ -discounted decision problem, whose reward has been shaped by the solution of the decision problem at the previous iteration. Unlike the *biased* solution obtained by solving the decision problem with a smaller value of γ (discussed above), by iteratively solving decision problems with a shorter $\gamma\kappa$ horizon, the κ -PI and κ -VI algorithms could converge to an optimal policy of the original decision problem.

In this paper, we derive model-free RL algorithms based on the κ -greedy formulation of multi-step greedy policies. As mentioned earlier, the main component of this formulation is (approximately) solving a surrogate decision problem with a shaped reward and a smaller discount factor. Our algorithms build on κ -PI and κ -VI, and solve the surrogate decision problem with the popular deep RL algorithms: Deep Q-Network (DQN) (Mnih et al., 2015) and Trust Region Policy Optimization (TRPO) (Schulman et al., 2015). We call the resulting algorithms κ -PI-DQN, κ -VI-DQN, κ -PI-TRPO, and κ -VI-TRPO, and empirically evaluate and compare them with DQN, TRPO, and Generalized Advantage Estimation (GAE) (Schulman et al., 2016) on Atari (Bellemare et al., 2013) and MuJoCo (Todorov et al., 2012) benchmarks. Our results indicate that for the right range of κ , our algorithms outperform DQN and TRPO. This suggests that the performance of these two deep RL algorithms can be

^{*}Equal contribution ¹Facebook AI Research, Menlo Park, USA

²Technion, Haifa, Israel ³Google Research, Mountain View, USA.
Correspondence to: Manan Tomar <manan.tomar@gmail.com>.

improved by using them as a solver within the multi-step greedy PI and VI schemes.

Moreover, our results indicate that the success of our algorithms depends on a number of non-trivial design choices. In particular, we identify the importance of a hyper-parameter that controls the extent to which the surrogate decision problem is solved, and use the theory of multi-step greedy DP to derive a recipe for setting this parameter. We show the advantage of using *hard* over *soft* updates, verifying the theory in Efroni et al. (2018b, Thm. 1). By hard and soft update, we refer to fully solving the surrogate MDP in a model-free manner and then evaluating the resulting policy (policy improvement and evaluation steps are separated) vs. changing the policy at each iteration (policy improvement and evaluation steps are concurrent – each improvement is followed by an evaluation).

We also establish a connection between our multi-step greedy algorithms and GAE. In particular, we show that our κ -PI-TRPO algorithm coincides with GAE and we can obtain GAE by minor modifications to κ -PI-TRPO. Finally, we show the advantage of using our multi-step greedy algorithms over lowering the discount factor in DQN (value-based) and TRPO (policy-based) algorithms. Our results indicate that while lowering the discount factor is detrimental to performance, our multi-step greedy algorithms indeed improve over DQN and TRPO.

2. Preliminaries

In this paper, we assume that the agent’s interaction with the environment is modeled as a discrete time γ -discounted Markov Decision Process (MDP), defined by $\mathcal{M}_\gamma = (\mathcal{S}, \mathcal{A}, P, R, \gamma, \mu)$, where \mathcal{S} and \mathcal{A} are the state and action spaces; $P \equiv P(s'|s, a)$ is the transition kernel; $R \equiv r(s, a)$ is the reward function with the maximum value of R_{\max} ; $\gamma \in (0, 1)$ is the discount factor; and μ is the initial state distribution. Let $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$ be a stationary Markovian policy, where $\mathcal{P}(\mathcal{A})$ is a probability distribution on the set \mathcal{A} . The value function of a policy π at any state $s \in \mathcal{S}$ is defined as $V^\pi(s) \equiv \mathbb{E}[\sum_{t \geq 0} \gamma^t r(s_t, a_t) | s_0 = s, \pi]$, where the expectation is over all the randomness in the policy, dynamics, and rewards. Similarly, the action-value function of π is defined as $Q^\pi(s, a) = \mathbb{E}[\sum_{t \geq 0} \gamma^t r(s_t, a_t) | s_0 = s, a_0 = a, \pi]$. Since the rewards are bounded by R_{\max} , both V and Q functions have the maximum value of $V_{\max} = R_{\max}/(1 - \gamma)$. An optimal policy π^* is the policy with maximum value at every state. We call the value of π^* the optimal value, and define it as $V^*(s) = \max_\pi V^\pi(s)$, $\forall s \in \mathcal{S}$. We denote by $Q^*(s, a)$, the state-action value of π^* , and remind that the following relation holds $V^*(s) = \max_a Q^*(s, a)$, for all s . The algorithms by which we solve an MDP (obtain an optimal policy) are mainly based on two popular DP algorithms: Policy Iteration (PI) and Value Iteration (VI). While

VI relies on iteratively computing the optimal Bellman operator \mathcal{T} applied to the current value function V (Eq. 1), PI relies on (iteratively) calculating a 1-step greedy policy $\pi_{1\text{-step}}$ w.r.t. to the value function of the current policy V (Eq. 2): for all $s \in \mathcal{S}$, we have

$$(\mathcal{T}V)(s) = \max_{a \in \mathcal{A}} \mathbb{E}[r(s_0, a) + \gamma V(s_1) | s_0 = s], \quad (1)$$

$$\pi_{1\text{-step}}(s) \in \arg \max_{a \in \mathcal{A}} \mathbb{E}[r(s_0, a) + \gamma V(s_1) | s_0 = s]. \quad (2)$$

It is known that \mathcal{T} is a γ -contraction w.r.t. the max-norm and its unique fixed-point is V^* , and the 1-step greedy policy w.r.t. V^* is an optimal policy π^* . In practice, the state space is often large, and thus, we can only approximately compute Eqs. 1 and 2, which results in approximate PI (API) and VI (AVI) algorithms. These approximation errors then propagate through the iterations of the API and AVI algorithms. However, it has been shown that this (propagated) error can be controlled (Munos, 2003; 2005; Farahmand et al., 2010), and after N steps, the algorithms approximately converge to a solution π_N , whose difference with the optimal value is bounded (see e.g., Scherrer, 2014 for API):

$$\eta(\pi^*) - \eta(\pi_N) \leq C\delta/(1 - \gamma)^2 + \gamma^N V_{\max}. \quad (3)$$

In Eq. 3, the scalar $\eta(\pi) = \mathbb{E}_{s \sim \mu}[V^\pi(s)]$ is the expected value function at the initial state,¹ δ represents the per-iteration error, and C upper-bounds the mismatch between the sampling distribution and the distribution according to which the final value function is evaluated (μ in Eq. 3), depending heavily on the dynamics. Finally, the second term on the RHS of Eq. 3 is the error due to initial values of policy/value and decays with the number of iterations N .

3. κ -Greedy Policy: κ -PI and κ -VI Algorithms

The optimal Bellman operator \mathcal{T} (Eq. 1) and 1-step greedy policy $\pi_{1\text{-step}}$ (Eq. 2) can be generalized to their multi-step versions. The most straightforward form of this generalization is realized by replacing \mathcal{T} and $\pi_{1\text{-step}}$ with h -optimal Bellman operator and h -step greedy policy (i.e., a lookahead of horizon h), respectively. This is done by substituting the 1-step return in Eqs. 1 and 2, $r(s_0, a) + \gamma V(s_1)$, with the h -step return, $\sum_{t=0}^{h-1} r(s_t, a_t) + \gamma^h V(s_h)$, and computing the maximum over actions a_0, \dots, a_{h-1} , instead of just a_0 (Bertsekas & Tsitsiklis, 1996). Efroni et al. (2018a) proposed an alternative form for the multi-step optimal Bellman operator and greedy policy, called κ -optimal Bellman oper-

¹Note that the LHS of Eq. 3 is the ℓ_1 -norm of $(V^{\pi^*} - V^{\pi_N})$ w.r.t. the initial state distribution μ .

ator, \mathcal{T}_κ , and κ -greedy policy, π_κ , for $\kappa \in [0, 1]$, i.e.,

$$(\mathcal{T}_\kappa V)(s) = \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} (\gamma \kappa)^t r_t(\kappa, V) \mid s_0 = s, \pi \right], \quad (4)$$

$$\pi_\kappa(s) \in \arg \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} (\gamma \kappa)^t r_t(\kappa, V) \mid s_0 = s, \pi \right], \quad (5)$$

for all $s \in \mathcal{S}$. In Eqs. 4 and 5, the *shaped reward* $r_t(\kappa, V)$ w.r.t. the value function V is defined as

$$r_t(\kappa, V) \equiv r_t + \gamma(1 - \kappa)V(s_{t+1}). \quad (6)$$

It can be shown that the κ -greedy policy w.r.t. the value function V is the optimal policy w.r.t. a κ -weighted geometric average of all future h -step returns (from $h = 0$ to ∞). This can be interpreted as TD(λ) (Sutton & Barto, 2018) for policy improvement (see Efroni et al., 2018a, Sec. 6). The important difference is that TD(λ) is used for policy evaluation and not for policy improvement.

It is easy to see that solving Eqs. 4 and 5 is equivalent to solving a surrogate $\gamma\kappa$ -discounted MDP with the shaped reward $r_t(\kappa, V)$, which we denote by $\mathcal{M}_{\gamma\kappa}(V)$ throughout the paper. The optimal value and policy of the surrogate MDP $\mathcal{M}_{\gamma\kappa}(V)$ are $\mathcal{T}_\kappa V$ and the κ -greedy policy π_κ , respectively. Using the notions of κ -optimal Bellman operator, \mathcal{T}_κ , and κ -greedy policy, π_κ , Efroni et al. (2018a) derived κ -PI and κ -VI algorithms, whose pseudocodes are shown in Algorithms 1 and 2. κ -PI iteratively (i) evaluates the value V^{π_i} of the current policy π_i , and (ii) sets the new policy, π_{i+1} , to the κ -greedy policy w.r.t. the value of the current policy V^{π_i} , by solving Eq. 5. On the other hand, κ -VI repeatedly applies the \mathcal{T}_κ operator to the current value function V_i (solves Eq. 4) to obtain the next value function, V_{i+1} , and returns the κ -greedy policy w.r.t. the final value V_{N_κ} . Note that for $\kappa = 0$, the κ -optimal Bellman operator and κ -greedy policy are equivalent to their 1-step counterparts, defined by Eqs. 1 and 2, which indicates that κ -PI and κ -VI are generalizations of the seminal PI and VI algorithms.

Algorithm 1 κ -Policy Iteration

```

1: Initialize:  $\kappa \in [0, 1]$ ,  $\pi_0$ ,  $N_\kappa$ 
2: for  $i = 0, 1, \dots, N_\kappa - 1$  do
3:    $V^{\pi_i} = \mathbb{E}[\sum_{t \geq 0} \gamma^t r_t \mid \pi_i]$ 
4:    $\pi_{i+1} \leftarrow \arg \max_{\pi} \mathbb{E}[\sum_{t \geq 0} (\gamma \kappa)^t r_t(\kappa, V^{\pi_i}) \mid \pi]$ 
5: end for
6: Return  $\pi_{N_\kappa}$ 
    
```

Algorithm 2 κ -Value Iteration

```

1: Initialize:  $\kappa \in [0, 1]$ ,  $V_0$ ,  $N_\kappa$ 
2: for  $i = 0, 1, \dots, N_\kappa - 1$  do
3:    $V_{i+1} = \max_{\pi} \mathbb{E}[\sum_{t \geq 0} (\gamma \kappa)^t r_t(\kappa, V_i) \mid \pi]$ 
4: end for
5:  $\pi_{N_\kappa} \leftarrow \arg \max_{\pi} \mathbb{E}[\sum_{t \geq 0} (\gamma \kappa)^t r_t(\kappa, V_{N_\kappa}) \mid \pi]$ 
6: Return  $\pi_{N_\kappa}$ 
    
```

It has been shown that both PI and VI converge to the optimal value with an exponential rate that depends on the discount factor γ , i.e., $\|V^* - V^{\pi_N}\|_\infty \leq O(\gamma^N)$ (see e.g., Bertsekas & Tsitsiklis, 1996). Analogously, Efroni et al. (2018a) showed that κ -PI and κ -VI converge with a faster exponential rate $\xi_\kappa = \frac{\gamma(1-\kappa)}{1-\gamma\kappa} \leq \gamma$, i.e., $\|V^* - V^{\pi_{N_\kappa}}\|_\infty \leq O(\xi_\kappa^{N_\kappa})$, with the cost that each iteration of these algorithms is computationally more expensive than that of PI and VI. Finally, we state the following property of κ -PI and κ -greedy policies that we use in our κ -PI and κ -VI based RL algorithms described in Section 4:

Asymptotic performance depends on κ . Efroni et al. (2018b, Thm. 5) proved the following bound on the performance of κ -PI that is similar to the one in Eq. 3 for API:

$$\eta(\pi^*) - \eta(\pi_{N_\kappa}) \leq \underbrace{C_\kappa \delta_\kappa / (1 - \gamma)^2}_{\text{Asymptotic Term}} + \underbrace{\xi_\kappa^{N_\kappa} V_{\max}}_{\text{Decaying Term}}, \quad (7)$$

where δ_κ and C_κ are quantities similar to δ and C in Eq. 3. Note that while the second term on the RHS of Eq. 7 decays with N_κ , the first one is independent of N_κ .

4. κ -PI and κ -VI based RL Algorithms

As described in Section 3, implementing κ -PI and κ -VI requires iteratively solving a $\gamma\kappa$ -discounted surrogate MDP with a shaped reward. If a model of the problem is given, the surrogate MDP can be solved using a DP algorithm (see Efroni et al., 2018a, Sec. 7). When a model is not available, we should approximately solve the surrogate MDP using a model-free RL algorithm. In this paper, we focus on the latter case and propose RL algorithms inspired by κ -PI and κ -VI. In our algorithms, we use model-free RL algorithms DQN (Mnih et al., 2015) and TRPO (Schulman et al., 2015) (for discrete and continuous action problems, respectively) as subroutines for estimating a κ -greedy policy (Line 4 in Alg. 1, κ -PI, and Line 5 in Alg. 2, κ -VI) and an optimal value of the surrogate MDP (Line 3 in Alg. 2, κ -VI). We refer to the resulting algorithms as κ -PI-DQN, κ -VI-DQN, κ -PI-TRPO, and κ -VI-TRPO.

In order to have an efficient implementation of our κ -PI and κ -VI based algorithms, the main question to answer is how a fixed number of samples T should be allocated to different parts of the κ -PI and κ -VI algorithms? More precisely, how shall we set $N_\kappa \in \mathbb{N}$, i.e., the total number of iterations of our algorithms, and determine the number of samples to solve the surrogate MDP at each iteration? To answer these questions, we devise a heuristic approach based on the theory of κ -PI and κ -VI algorithms, and in particular Eq. 7. Since N_κ only appears explicitly in the second term on the RHS of Eq. 7, an appropriate choice of N_κ is such that $C_\kappa \delta_\kappa / (1 - \gamma)^2 \simeq \xi_\kappa^{N_\kappa} V_{\max}$. Note that setting N_κ to a higher value would *not* significantly improve the performance, because the asymptotic term in Eq. 7 is independent

of N_κ . In practice, since δ_κ and C_κ are unknown, we set N_κ to satisfy the following equality:

$$\xi_\kappa^{N_\kappa} = C_{\text{FA}}, \quad (8)$$

where C_{FA} is a hyper-parameter that depends on the *final-accuracy* we aim for. For example, if our goal is to obtain 90% accuracy, we would set $C_{\text{FA}} = 0.1$, which results in $N_{\kappa=0.99} \simeq 4$ and $N_{\kappa=0.5} \simeq 115$, for $\gamma = 0.99$. Our experimental results in Section 5 suggest that this approach leads to a reasonable choice for the total number of iterations N_κ . It is important to note the following facts: **1)** as we increase κ , we expect less iterations are needed for κ -PI and κ -VI to converge to a good policy, and **2)** the *effective horizon*² of the surrogate MDP that κ -PI and κ -VI solve at each iteration increases with κ .

Lastly, we need to determine the number of samples for each iteration of our κ -PI and κ -VI based algorithms. We allocate equal number of samples per iteration, denoted by T_κ . Since the total number of samples, T , is known beforehand, we set the number of samples per iteration to

$$T_\kappa = T/N_\kappa. \quad (9)$$

In the rest of the paper, we first derive our DQN-based and TRPO-based algorithms in Sections 4.1 and 4.2. It is important to note that for $\kappa = 1$, our algorithms are reduced to DQN and TRPO. We then conduct a set of experiments with our algorithms in Sections 5.1 and 5.2 in which we carefully study the effect of κ and N_κ (or equivalently the hyper-parameter C_{FA} , defined by Eq. 8) on their performance.

4.1. κ -PI-DQN and κ -VI-DQN Algorithms

Algorithm 3 presents the pseudo-code of κ -PI-DQN. Due to space constraints, we report the detailed pseudo-code in Appendix A.1 (Alg. 5). We use four neural networks in this algorithm, two to represent the Q -function of the original MDP (with discount factor γ and reward r), Q_ϕ (Q -network) and Q'_ϕ (target network), and two for the Q -function of the surrogate MDP, Q_θ (Q -network) and Q'_θ (target network). In the *policy improvement step*, we use DQN to solve the $\gamma\kappa$ -discounted surrogate MDP with the shaped reward $r_j(\kappa, V_\phi) = r_j + \gamma(1 - \kappa)V_\phi(s_{j+1})$, i.e., $\mathcal{M}_{\gamma\kappa}(V_\phi)$, where $V_\phi \simeq V^{\pi_{i-1}}$ and is computed as $V_\phi(s) = Q_\phi(s, \arg \max_a Q'_\theta(s, a))$. The output of DQN is (approximately) the optimal Q -function of $\mathcal{M}_{\gamma\kappa}(V_\phi)$, and thus, the new policy π_i , which is the (approximate) κ -greedy policy w.r.t. V_ϕ is equal to $\pi_i(\cdot) = \arg \max_a Q'_\theta(\cdot, a)$. In the *policy evaluation step*, we use off-policy TD(0) to evaluate the Q -function of the current policy π_i , i.e., $Q_\phi \simeq Q^{\pi_i}$. Although what is needed is an estimate of the value function

²The effective horizon of a $\gamma\kappa$ -discounted MDP is $1/(1 - \gamma\kappa)$.

Algorithm 3 κ -PI-DQN

```

1: Initialize replay buffer  $\mathcal{D}$ ;  $Q$ -networks  $Q_\theta, Q_\phi$ ; target net-
   works  $Q'_\theta, Q'_\phi$ ;
2: for  $i = 0, \dots, N_\kappa - 1$  do
3:   # Policy Improvement
4:   for  $t = 1, \dots, T_\kappa$  do
5:     Act by an  $\epsilon$ -greedy policy w.r.t.  $Q_\theta(s_t, a)$ , observe
        $r_t, s_{t+1}$ , and store  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ ;
6:     Sample a batch  $\{(s_j, a_j, r_j, s_{j+1})\}_{j=1}^N$  from  $\mathcal{D}$ ;
7:     Update  $\theta$  using DQN with
        $\{(s_j, a_j, r_j(\kappa, V_\phi), s_{j+1})\}_{j=1}^N$ , where
8:      $V_\phi(s_{j+1}) = Q_\phi(s_{j+1}, \pi_{i-1}(s_{j+1}))$  and
9:      $\pi_{i-1}(\cdot) \in \arg \max_a Q'_\theta(\cdot, a)$ ;
10:    Copy  $\theta$  to  $\theta'$  occasionally ( $\theta' \leftarrow \theta$ );
11:   end for
12:   # Policy Evaluation of  $\pi_i(s) \in \arg \max_a Q'_\theta(s, a)$ 
13:   for  $t = 1, \dots, T_\kappa$  do
14:     Sample a batch  $\{(s_j, a_j, r_j, s_{j+1})\}_{j=1}^N$  from  $\mathcal{D}$ ;
15:     Update  $\phi$  using this data and off-policy TD(0) to estimate
       the  $Q$ -function of the current policy  $\pi_i$ ;
16:     Copy  $\phi$  to  $\phi'$  occasionally ( $\phi' \leftarrow \phi$ );
17:   end for
18: end for

```

of the current policy, $V_\phi \simeq V^{\pi_i}$, we chose to evaluate its Q -function, because the available data (the transitions stored in the replay buffer) is off-policy, and unlike the value function, the Q -function of a fixed policy can be easily evaluated with this type of data using off-policy TD(0).

Remark 1. In the policy improvement phase, π_{i-1} is computed as the $\arg \max_a Q'_\theta(\cdot, a)$ (Line 10), a quantity that is not constant and is (slowly) changing during this phase. Addressing this issue requires using an additional target network that is set to Q_θ only at the end of each improvement step and its $\arg \max$ is used to compute π_{i-1} throughout the improvement step of the next iteration. We tested using this additional network in our experiments, but it did not improve the performance, and thus, we decided to report the algorithm without it.

We report the pseudo-code of κ -VI-DQN in Appendix A.1 (Alg. 6). Note that κ -VI simply repeats $V \leftarrow \mathcal{T}_\kappa V$ and computes $\mathcal{T}_\kappa V$, which is the optimal value of the surrogate MDP $\mathcal{M}_{\gamma\kappa}(V)$. In κ -VI-DQN, we repeatedly solve $\mathcal{M}_{\gamma\kappa}(V)$ by DQN and use its (approximately) optimal Q -function to shape the reward of the next iteration. The algorithm uses three neural networks, two to solve the surrogate MDP by DQN, Q_θ (Q -network) and Q'_θ (target network), and one to store its optimal Q -function to use it for the shaped reward in the next iteration, Q_ϕ . Let $Q_{\gamma\kappa, V}^*$ and $V_{\gamma\kappa, V}^*$ be the optimal Q and V functions of $\mathcal{M}_{\gamma\kappa}(V)$. Then, we have $\max_a Q_{\gamma\kappa, V}^*(s, a) = V_{\gamma\kappa, V}^*(s) = (\mathcal{T}_\kappa V)(s)$, where the first equality is by definition (Sec. 2) and the second one holds since $\mathcal{T}_\kappa V$ is the optimal value of $\mathcal{M}_{\gamma\kappa}(V)$ (Sec. 3). Therefore, in κ -VI-DQN, we shape the reward at each iteration by $\max_a Q_\phi(s, a)$, where Q_ϕ is the output of the DQN

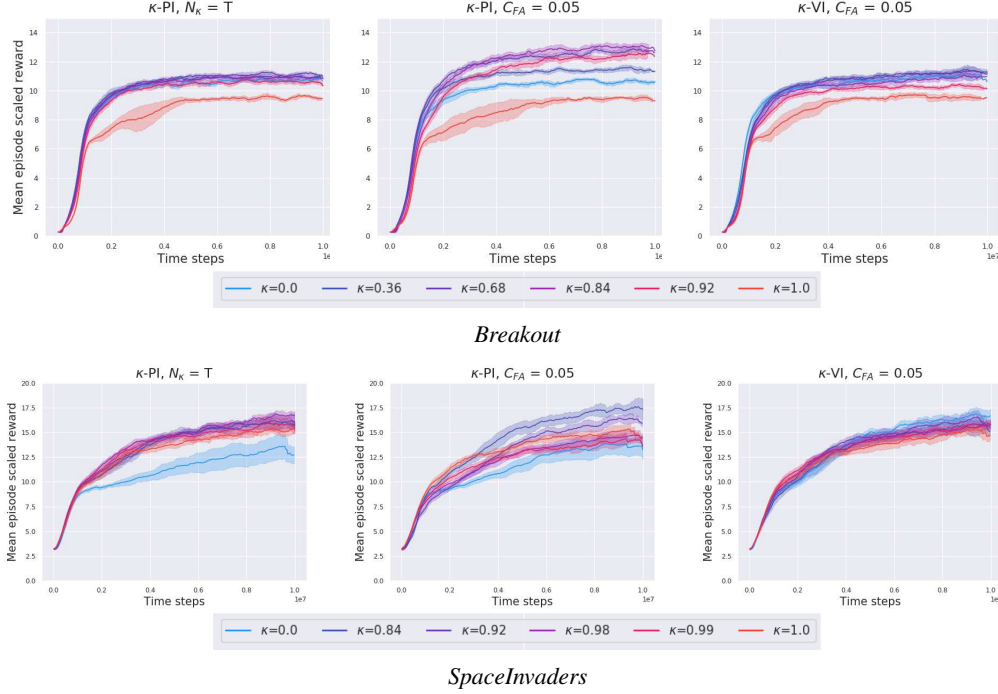


Figure 1: Training performance of the ‘naive’ baseline $N_\kappa = T$ and κ -PI-DQN, κ -VI-DQN for $C_{FA} = 0.05$ on Breakout (top) and SpaceInvaders (bottom). See Appendix A.2 for performance w.r.t. different C_{FA} values.

Algorithm 4 κ -PI-TRPO

```

1: Initialize  $V$ -networks  $V_\theta$  and  $V_\phi$ ; policy network  $\pi_\psi$ ;
2: for  $i = 0, \dots, N_\kappa - 1$  do
3:   for  $t = 1, \dots, T_\kappa$  do
4:     Simulate the current policy  $\pi_\psi$  for  $M$  steps and calculate
       the following two returns for all steps  $j$ :
        $R_j(\kappa, V_\phi) = \sum_{t=j}^M (\gamma \kappa)^{t-j} r_t(\kappa, V_\phi)$  and
        $\rho_j = \sum_{t=j}^M \gamma^{t-j} r_t$ ;
5:     Update  $\theta$  by minimizing the batch loss function:
        $\mathcal{L}_{V_\theta} = \frac{1}{N} \sum_{j=1}^N (V_\theta(s_j) - R_j(\kappa, V_\phi))^2$ ;
6:     # Policy Improvement
7:     Update  $\psi$  using TRPO and the batch
        $\{(R_j(\kappa, V_\phi), V_\theta(s_j))\}_{j=1}^N$ ;
8:   end for
9:   # Policy Evaluation
10:  Update  $\phi$  by minimizing the batch loss function:
        $\mathcal{L}_{V_\phi} = \frac{1}{N} \sum_{j=1}^N (V_\phi(s_j) - \rho_j)^2$ ;
11: end for
    
```

from the previous iteration, i.e., $\max_a Q_\phi(s, a) \simeq \mathcal{T}_\kappa V_{i-1}$.

4.2. κ -PI-TRPO and κ -VI-TRPO Algorithms

Algorithm 4 presents the pseudo-code of κ -PI-TRPO (a detailed pseudo-code is reported in Appendix B.1, Alg. 7). Recall that TRPO iteratively updates the current policy using its return and an estimate of its value function. At each iteration i of κ -PI-TRPO: **1)** we use the estimate of the

value of the current policy $V_\phi \simeq V^{\pi_{i-1}}$ to calculate the return $R(\kappa, V_\phi)$ and the estimate of the value function V_θ of the surrogate MDP $\mathcal{M}_{\gamma\kappa}(V_\phi, \mathbf{2})$ we use $R(\kappa, V_\phi)$ and V_θ to compute the new policy π_i (TRPO style), and **3)** we estimate the value of the new policy $V_\phi \simeq V^{\pi_i}$ in the original MDP (with discount factor γ and reward r). The algorithm uses three neural networks, one for the value function of the original MDP, V_ϕ , one for the value function of the surrogate MDP, V_θ , and one for the policy, π_ψ .

We report the pseudo-code of κ -VI-TRPO in Appendix B.1, Alg. 8. As previously noted, κ -VI iteratively solves the surrogate MDP and uses its optimal value $\mathcal{T}_\kappa V_{i-1}$ to shape the reward of the surrogate MDP in the next iteration. In κ -VI-TRPO, we solve the surrogate MDP $\mathcal{M}_{\gamma\kappa}(V_{i-1} = V_\phi)$ with TRPO until its policy π_ψ converges to the optimal policy of $\mathcal{M}_{\gamma\kappa}(V_{i-1} = V_\phi)$ and its value V_θ converges to $\mathcal{T}_\kappa V_{i-1} = \mathcal{T}_\kappa V_\phi$. We then replace V_i with V_θ ($V_i = V_\phi \leftarrow V_\theta$) and repeat this process.

5. Experimental Results

In our experiments, we specifically focus on answering the following questions:

1. Does the performance of DQN and TRPO improve when using them as κ -greedy solvers in κ -PI and κ -VI? Is there a performance tradeoff w.r.t. to κ ?

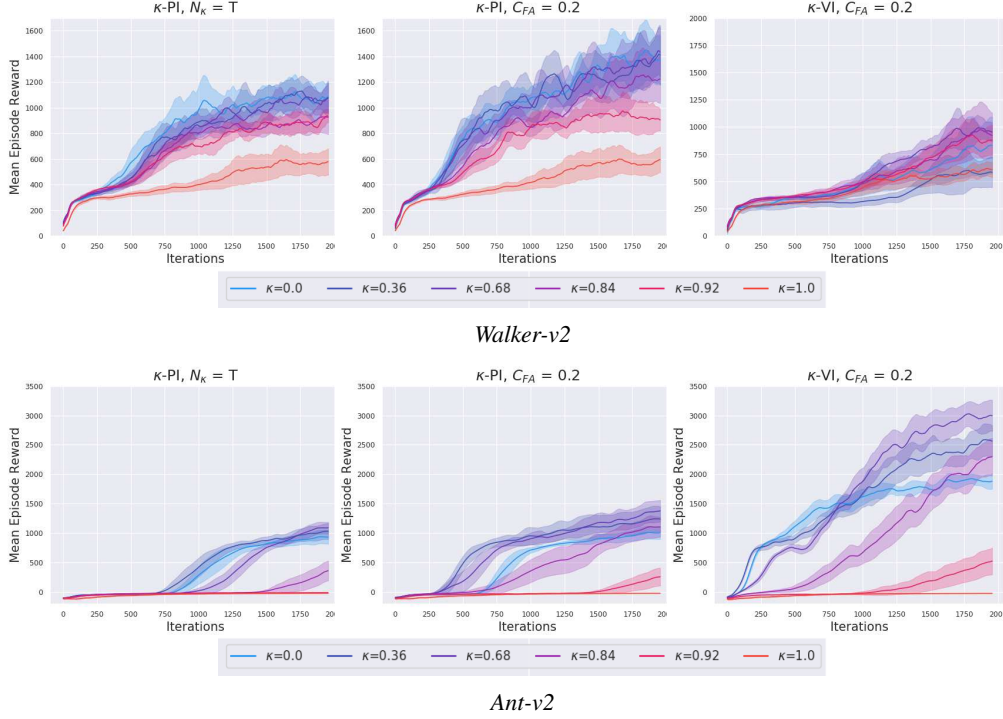


Figure 2: Training performance of the ‘naive’ baseline $N_\kappa = T$ and κ -PI-TRPO, κ -VI-TRPO for $C_{FA} = 0.2$ on Walker (top) and Ant (bottom). See Appendix B.2 for performance w.r.t. different C_{FA} values. Each iteration corresponds to roughly 1000 environment samples, and thus, the total number of training samples is 2 millions.

- Following κ -PI and κ -VI, our DQN and TRPO implementations of these algorithms devote a significant T_κ number of samples to each iteration. Is this needed or a ‘naive’ choice of $T_\kappa = 1$, or equivalently $N_\kappa = T$, works just well for all values of κ ?

We choose to test our κ -DQN and κ -TRPO algorithms on the Atari and MuJoCo benchmarks, respectively. Both of these algorithms use standard setups, including the use of the Adam optimizer for performing gradient descent, a discount factor of 0.99 across all tasks, target Q value networks in the case of κ -DQN and an entropy regularizer with a coefficient of 0.01 in the case of κ -TRPO. We choose to run our experiments for multiple values of κ between zero and one, which roughly follow a logarithmic scale. Note that just by using the definition of κ -greedy algorithms, we reduce to the base cases of DQN and TRPO when setting $\kappa = 1.0$ for all experiments. This forms as one of the two baselines we consider in this work. The second baseline essentially refers to using our κ -greedy algorithms for a fixed $N_\kappa = T$ value. Thus, independent of the value of κ , the surrogate MDP is solved for a single time-step per each iteration. Below, we describe the experiments and results in further detail. The implementation details for the κ -DQN and κ -TRPO cases are provided in Appendix A, Table 3 and Appendix B, Table 4, respectively.

5.1. κ -PI-DQN and κ -VI-DQN Experiments

In this section, we empirically analyze the performance of the κ -PI-DQN and κ -VI-DQN algorithms on the Atari domains: Breakout, SpaceInvaders, Seaquest, Enduro, BeamRider, and Qbert (Bellemare et al., 2013). We start by performing an ablation test on three values of hyper-parameter $C_{FA} = \{0.001, 0.05, 0.2\}$ on the Breakout domain. The value of C_{FA} sets the number of samples per iteration T_κ (Eq. 8) and the total number of iterations N_κ (Eq. 9). The total number of samples is set to $T \simeq 10^6$. This value represents the number of samples after which our DQN-based algorithms approximately converge. For each value of C_{FA} , we test κ -PI-DQN and κ -VI-DQN for several κ values. In both algorithms, the best performance was obtained with $C_{FA} = 0.05$. Therefore, C_{FA} is set to 0.05 for all our experiments with other Atari domains.

Figure 1 shows the training performance of κ -PI-DQN and κ -VI-DQN on Breakout and SpaceInvaders for the best value of $C_{FA} = 0.05$, as well as for the ‘naive’ baseline $T_\kappa = 1$, or equivalently $N_\kappa = T$. The results on Breakout for the other values of C_{FA} and the training plots for all other Atari domains (with $C_{FA} = 0.05$) are reported in Appendices A.2 and A.3, respectively. Table 1 shows the final training performance of κ -PI-DQN and κ -VI-DQN

Multi-step Greedy Reinforcement Learning Algorithms

| Domain | Alg. | κ_{best} | $\kappa = 0$ | DQN, $\kappa = 1$ | $N_\kappa = T$, κ_{best} (κ -PI) |
|-----------|--------------|--|--------------------|-------------------|---|
| Breakout | κ -PI | 223 (± 7), $\kappa=0.84$ | 154(± 3) | 134(± 4) | 170(± 2), $\kappa=0.68$ |
| | κ -VI | 181(± 7), $\kappa=0.68$ | 174(± 5) | | |
| SpaceInv. | κ -PI | 755 (± 23), $\kappa=0.84$ | 613(± 20) | 656(± 17) | 700(± 21), $\kappa=0.98$ |
| | κ -VI | 712(± 25), $\kappa=0.92$ | 687(± 32) | | |
| Seaquest | κ -PI | 5159 (± 292), $\kappa=0.92$ | 2612(± 238) | 3099(± 191) | 3897(± 218), $\kappa=0.68$ |
| | κ -VI | 3253(± 402), $\kappa=0.84$ | 2680(± 382) | | |
| Enduro | κ -PI | 533 (± 12), $\kappa=0.84$ | 478(± 10) | 224(± 110) | 535 (± 13), $\kappa=0.68$ |
| | κ -VI | 486(± 23), $\kappa=0.84$ | 443 (± 90) | | |
| BeamRider | κ -PI | 3849(± 110), $\kappa=1.0$ | 3103(± 279) | 3849(± 110) | 3849(± 110), $\kappa=1.0$ |
| | κ -VI | 4277 (± 269), $\kappa=0.84$ | 3714 (± 437) | | |
| Qbert | κ -PI | 8157 (± 265), $\kappa=0.84$ | 6719(± 520) | 7258(± 385) | 7968 (± 218), $\kappa=0.98$ |
| | κ -VI | 8060 (± 158), $\kappa=0.84$ | 7563 (± 398) | | |

Table 1: The final training performance of κ -PI-DQN and κ -VI-DQN on the Atari domains, for the hyper-parameter $C_{FA} = 0.05$. The values are reported for a 95% confidence interval across 10 random runs (*empirical mean* $\pm 1.96 \times$ *empirical standard deviation* / $\sqrt{n} = 10$). The best scores are in bold and multiple bold values for a domain denote an insignificant statistical difference between them.

| Domain | Alg. | κ_{best} | $\kappa = 0$ | TRPO, $\kappa = 1$ | $N_\kappa = T$, κ_{best} (κ -PI) | GAE, λ_{best} |
|---------------|--------------|--|---------------------|-------------------------|---|---|
| Walker | κ -PI | 1438 (± 188), $\kappa=0.68$ | 1371 (± 192) | 594 (± 99) | 1082(± 110), $\kappa=0.0$ | 1601 (± 190), $\lambda=0.0$ |
| | κ -VI | 954(± 88), $\kappa=0.68$ | 830(± 262) | | | |
| Ant | κ -PI | 1377(± 183), $\kappa=0.68$ | 1006(± 106) | -19(± 1) | 1090(± 99), $\kappa=0.68$ | 1094(± 139), $\lambda=0.0$ |
| | κ -VI | 2998 (± 264), $\kappa=0.68$ | 1879(± 128) | | | |
| HalfCheetah | κ -PI | 1334 (± 151), $\kappa=0.36$ | 907(± 176) | -18(± 87) | 1195(± 218), $\kappa=0.36$ | 1322 (± 213), $\lambda=0.36$ |
| | κ -VI | 1447 (± 346), $\kappa=0.36$ | 1072(± 30) | | | |
| HumanoidStand | κ -PI | 72604 (± 1219), $\kappa=0.99$ | 52936(± 1529) | 68143(± 1031) | 71331 (± 1149), $\kappa=0.98$ | 71932 (± 2122), $\lambda=0.98$ |
| | κ -VI | 72821 (± 908), $\kappa=0.99$ | 51148(± 1377) | | | |
| Swimmer | κ -PI | 107 (± 12), $\kappa=1.0$ | 42(± 3) | 107 (± 12) | 107 (± 12), $\kappa=1.0$ | 103 (± 13), $\lambda = 1.0$ |
| | κ -VI | 114 (± 15), $\kappa=1.0$ | 46(± 1) | | | |
| Hopper | κ -PI | 1486 (± 324), $\kappa=0.68$ | 1012(± 263) | 1142(± 141) | 1434 (± 129), $\kappa=0.98$ | 1600 (± 134), $\lambda = 0.84$ |
| | κ -VI | 1069(± 76), $\kappa=0.92$ | 531(± 125) | | | |

Table 2: The final training performance of κ -PI-TRPO and κ -VI-TRPO on the MuJoCo domains, for the hyper-parameter $C_{FA} = 0.2$. The values are reported for a 95% confidence interval across 10 random runs (*empirical mean* $\pm 1.96 \times$ *empirical standard deviation* / $\sqrt{n} = 10$). The best scores are in bold and multiple bold values for a domain denote an insignificant statistical difference between them.

on the Atari domains with $C_{FA} = 0.05$. Note that the scores reported in Table 1 are the actual returns on the Atari domains, while the vertical axis in the plots of Figure 1 corresponds to a scaled return. We plot the scaled return, since this way it can be easier to reproduce our results using the OpenAI Baselines codebase (Hill et al., 2018).

Our results exhibit that both κ -PI-DQN and κ -VI-DQN improve the performance of DQN ($\kappa = 1$). Moreover, they show that setting $N_\kappa = T$ leads to a clear degradation of the final training performance on all of the domains except for Enduro, which gets to approximately the same score. Although the performance degrades, the results for $N_\kappa = T$ are still better than for DQN.

5.2. κ -PI-TRPO and κ -VI-TRPO Experiments

In this section, we empirically analyze the performance of the κ -PI-TRPO and κ -VI-TRPO algorithms on the MuJoCo (Todorov et al., 2012) based OpenAI Gym domains:

Walker2d-v2, Ant-v2, HalfCheetah-v2, HumanoidStandup-v2, Swimmer-v2, and Hopper-v2 (Brockman et al., 2016). As in Section 5.1, we start by performing an ablation test on the parameter $C_{FA} = \{0.001, 0.05, 0.2\}$ on the Walker domain. We set the total number of iterations to 2000, with each iteration consisting 1000 samples. Thus, the total number of samples is $T \simeq 2 \times 10^6$. This is the number of samples after which our TRPO-based algorithms approximately converge. For each value of C_{FA} , we test κ -PI-TRPO and κ -VI-TRPO for several κ values. In both algorithms, the best performance was obtained with $C_{FA} = 0.2$, and thus, we set $C_{FA} = 0.2$ in our experiments with other MuJoCo domains.

Figure 2 shows the training performance of κ -PI-TRPO and κ -VI-TRPO on Walker and Ant domains for the best value of $C_{FA} = 0.2$, as well as for the ‘naive’ baseline $T_\kappa = 1$, or equivalently $N_\kappa = T$. The results on Walker for the other C_{FA} values and the other MuJoCo domains (with $C_{FA} = 0.2$) are reported in Appendices B.2 and B.3,

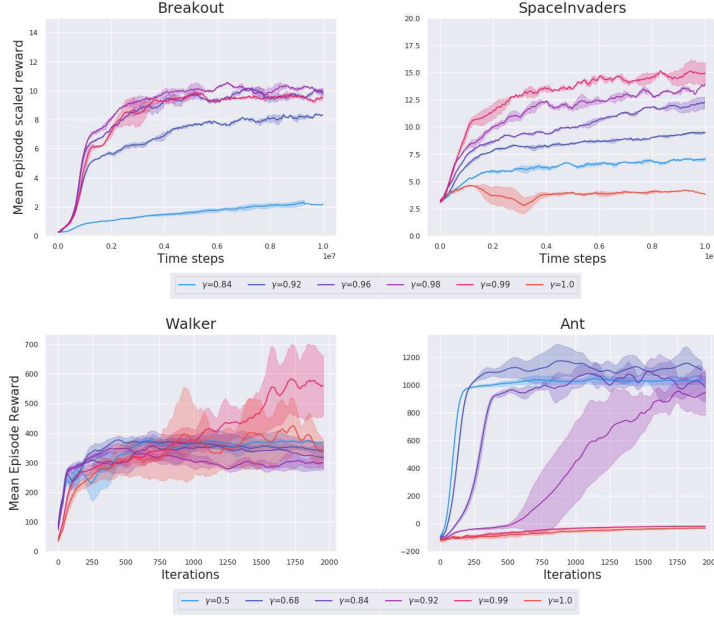


Figure 3: Lowering the discount factor γ for Atari domains Breakout and SpaceInvaders and MuJoCo domains Walker-v2 and Ant-v2

respectively. Table 2 shows the final training performance of κ -PI-TRPO and κ -VI-TRPO on the MuJoCo domains with $C_{FA} = 0.2$.

The results exhibit that both κ -PI-TRPO and κ -VI-TRPO yield better performance than TRPO ($\kappa = 1$). Furthermore, they show that the algorithms with $C_{FA} = 0.2$ perform better than with $N_\kappa = T$ for three out of six domains (Walker, Ant and HalfCheetah) and equally well for the remaining three (HumanoidStandup, Swimmer and Hopper).

6. Discussion and Related Work

Comparison with GAE: There is a close connection between the Generalized Advantage Estimation (GAE) algorithm (Schulman et al., 2016) and κ -PI. In GAE, the policy is updated by the following gradient:

$$\nabla_{\theta} \mathbb{E}_{s \sim \mu} [V^{\pi_{\theta}}(s)] = \mathbb{E}_{s \sim d_{\pi}^{\mu}} [\nabla_{\theta} \log \pi_{\theta}(s) \sum_t (\gamma \lambda)^t \delta_V],$$

$$\delta_V = r_t + \gamma V_{t+1} - V_t, \quad (10)$$

where d_{π}^{μ} is the occupancy measure of policy π . Eq. 10 can be interpreted as a gradient in a $\gamma \lambda$ -discounted MDP with shaped rewards δ_V , which we refer to as $\mathcal{M}_{\gamma \lambda}(\delta_V)$. As noted in Efroni et al. (2018a, Sec. 6), an optimal policy $\pi_{\gamma \lambda}^*$ of the MDP $\mathcal{M}_{\gamma \lambda}(\delta_V)$ is also optimal in $\mathcal{M}_{\gamma \kappa}(V)$ with $\kappa = \lambda$. This means that $\pi_{\gamma \lambda}^*$ is the κ -greedy policy w.r.t. V .

Comparing the two algorithms, we notice that GAE is conceptually similar to κ -PI-TRPO with $T_\kappa = 1$, or equivalently $N_\kappa = T$. In fact, in the case of $T_\kappa = 1$, we can

obtain a pseudo-code of GAE by removing Line 5 (no need to estimate the value of the surrogate MDP) and replacing $r_t(\kappa, V_{\phi})$ with the TD-error of the original MDP on Line 4 in Algorithm 4.

In Section 5.2, we compared the empirical performance of GAE with that of κ -PI-TRPO and κ -VI-TRPO (see Figure 2 and Table 2). The results show that κ -PI-TRPO and κ -VI-TRPO perform better than or on par with GAE. Moreover, we observe that in most domains, the performance of GAE is equivalent to that of κ -PI-TRPO with $N_\kappa = T$. This is in accordance with the description above connecting GAE to our naive baseline implementation. We report all GAE results in Appendix B.3.

Remark 2 (GAE Implementation). *In the OpenAI implementation of GAE, the value network is updated w.r.t. to the target $\sum_t (\gamma \lambda)^t r_t$, whereas in the GAE paper (Schulman et al., 2016), $\sum_t \gamma^t r_t$ is used as the target. We chose the latter form in our implementation of GAE to be in accord with the paper.*

Lowering Discount Factor in DQN and TRPO: To show the advantage of κ -PI and κ -VI based algorithms over simply lowering the discount factor γ , we test the performance of the “vanilla” DQN and TRPO algorithms with values of γ lower than the one previously used (i.e., $\gamma = 0.99$). As evident from Figure 3, only in the Ant domain, this approach results in an improved performance (for $\gamma = 0.68$). On the other hand, in the Ant domain, the performance of κ -PI-TRPO, and especially κ -VI-TRPO, surpasses that of TRPO

with the lower value of $\gamma = 0.68$. In Breakout, SpaceInvaders, and Walker, the performance of DQN and TRPO worsens or remains unchanged when we lower the discount factor (DQN and TRPO do not benefit from lowering γ), while our κ -PI and κ -VI based algorithms perform better with lowering the value of κ (note that our algorithms are reduced to DQN and TRPO for $\kappa = 1$).

Remark 3. While we observe better performance for smaller γ values in some MuJoCo domains, e.g., $\gamma = 0.68$ in Ant, lowering γ always results in inferior performance in the Atari domains. This is due to the fact that a number of MuJoCo domains, such as Ant, are inherently short-horizon decision problems, and thus, their performance does not degrade (even sometimes improves) with lowering the discount factor. On the other hand, the Atari problems are generally not short-horizon, and thus, their performance degrades with lowering γ .

7. Conclusion and Future Work

In this paper, we studied the use of multi-step greedy policies in model-free RL and showed that in most problems, the algorithms derived from this formulation achieve a better performance than their single-step counterparts. We adopted the κ -greedy formulation of multi-step greedy policies (Efroni et al., 2018a) and derived four model-free RL algorithms. The main component of the policy and value iteration algorithms derived from this formulation, κ -PI and κ -VI, is solving a surrogate decision problem with a shaped reward and a smaller discount factor. Our algorithms use popular deep RL algorithms, DQN and TRPO, to solve the surrogate decision problem, and thus, we refer to them as κ -PI-DQN, κ -VI-DQN, κ -PI-TRPO, and κ -VI-TRPO. We empirically evaluated our proposed algorithms and compared them with DQN, TRPO, and GAE on Atari and MuJoCo benchmarks. Our experiments show that for a large range of κ , our algorithms perform better than DQN and TRPO. Furthermore, we proposed a recipe to allocate the total sample budget to the evaluation and improvement phases of our algorithms, and empirically demonstrated the importance of this allocation. We also showed how GAE can be derived by minor modifications to κ -PI-TRPO, and thus, is a κ -greedy RL algorithm. Finally, we showed the advantage of multi-step greedy formulation over lowering the discount factor in DQN and TRPO. Our results indicate that while the performance of DQN and TRPO degrades with lowering the discount factor, our multi-step greedy algorithms improve over DQN and TRPO.

An interesting future direction would be to use other multi-step greedy formulations (Bertsekas & Tsitsiklis, 1996; Bertsekas, 2018; Efroni et al., 2018a; Sun et al., 2018; Shani et al., 2019) to derive model-free RL algorithms. Another direction is to use multi-step greedy in model-based RL

(e.g., Kumar et al., 2016; Talvitie, 2017; Luo et al., 2018; Janner et al., 2019) and solve the surrogate decision problem with an approximate model. We conjecture that in this case one may set κ – or more generally, the planning horizon – as a function of the quality of the approximate model: gradually increasing κ as the approximate model gets closer to the real one. We leave theoretical and empirical study of this problem for future work. Finally, we believe using adaptive κ would greatly improve the performance of our proposed algorithms. We leave verifying this and how κ should change as a function of errors in gradient and value estimation for future work.

References

- Baird, L. Residual algorithms: Reinforcement learning with function approximation. In *Proceedings of the Twelfth International Conference on Machine Learning*, pp. 30–37, 1995.
- Bellemare, M., Naddaf, Y., Veness, J., and Bowling, M. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- Bertsekas, D. Feature-based aggregation and deep reinforcement learning: A survey and some new implementations. *Journal of Automatica Sinica*, 6(1):1–31, 2018.
- Bertsekas, D. and Tsitsiklis, J. *Neuro-dynamic programming*, volume 5. 1996.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. OpenAI Gym. *Preprint arXiv:1606.01540*, 2016.
- Efroni, Y., Dalal, G., Scherrer, B., and Mannor, S. Beyond the one step greedy approach in reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning*, 2018a.
- Efroni, Y., Dalal, G., Scherrer, B., and Mannor, S. Multiple-step greedy policies in approximate and online reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 5238–5247, 2018b.
- Farahmand, A., Szepesvári, C., and Munos, R. Error propagation for approximate policy and value iteration. In *Advances in Neural Information Processing Systems*, pp. 568–576, 2010.
- Greensmith, E., Bartlett, P., and Baxter, J. Variance reduction techniques for gradient estimates in reinforcement learning. *Journal of Machine Learning Research*, 5:1471–1530, 2004.

- Hill, A., Raffin, A., Ernestus, M., Gleave, A., Kanervisto, A., Traore, R., Dhariwal, P., Hesse, C., Klimov, O., Nichol, A., Plappert, M., Radford, A., Schulman, J., Sidor, S., and Wu, Y. Stable baselines. GitHub repository <https://github.com/hill-a/stable-baselines>, 2018.
- Janner, M., Fu, J., Zhang, M., and Levine, S. When to trust your model: Model-based policy optimization. *arXiv preprint arXiv:1906.08253*, 2019.
- Jiang, N., Kulesza, A., Singh, S., and Lewis, R. The dependence of effective planning horizon on model accuracy. In *Proceedings of the International Conference on Autonomous Agents and Multiagent Systems*, pp. 1181–1189, 2015.
- Kumar, V., Todorov, E., and Levine, S. Optimal control with learned local models: Application to dexterous manipulation. In *International Conference on Robotics and Automation*, pp. 378–383, 2016.
- Levine, S., Finn, C., Darrell, T., and Abbeel, P. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- Lillicrap, T., Hunt, J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. *Preprint arXiv:1509.02971*, 2015.
- Lu, T., Schuurmans, D., and Boutilier, C. Non-delusional Q-learning and value iteration. In *Proceedings of the International Conference on Neural Information Processing Systems*, pp. 9971–9981, 2018.
- Luo, Y., Xu, H., Li, Y., Tian, Y., Darrell, T., and Ma, T. Algorithmic framework for model-based deep reinforcement learning with theoretical guarantees. *Preprint arXiv:1807.03858*, 2018.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M., Graves, A., Riedmiller, M., Fidjeland, A., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- Munos, R. Error bounds for approximate policy iteration. In *Proceedings of the International Conference on Machine Learning*, pp. 560–567, 2003.
- Munos, R. Error bounds for approximate value iteration. In *Proceedings of the National Conference on Artificial Intelligence*, pp. 1006–1011, 2005.
- Petrik, M. and Scherrer, B. Biasing approximate dynamic programming with a lower discount factor. In *Advances in neural information processing systems*, pp. 1265–1272, 2009.
- Scherrer, B. Approximate policy iteration schemes: a comparison. In *International Conference on Machine Learning*, pp. 1314–1322, 2014.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. Trust-region policy optimization. In *International conference on machine learning*, pp. 1889–1897, 2015.
- Schulman, J., Moritz, P., Levine, S., Jordan, M., and Abbeel, P. High-dimensional continuous control using generalized advantage estimation. In *Proceedings of the International Conference on Learning Representations*, 2016.
- Shani, L., Efroni, Y., and Mannor, S. Exploration conscious reinforcement learning revisited. In *International Conference on Machine Learning*, pp. 5680–5689, 2019.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., et al. Mastering the game of Go without human knowledge. *Nature*, 550(7676):354, 2017.
- Sun, W., Gordon, G., Boots, B., and Bagnell, J. Dual policy iteration. In *Advances in Neural Information Processing Systems*, pp. 7059–7069, 2018.
- Sutton, R. and Barto, A. *Reinforcement learning: An introduction*. 2018.
- Talvitie, E. Self-correcting models for model-based reinforcement learning. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- Thrun, S. and Schwartz, A. Issues in using function approximation for reinforcement learning. In *Proceedings of the Connectionist Models Summer School*, pp. 255–263, 1993.
- Todorov, E., Erez, T., and Tassa, Y. MuJoCo: A physics engine for model-based control. In *International Conference on Intelligent Robots and Systems*, pp. 5026–5033, 2012.
- Van Hasselt, H., Guez, A., and Silver, D. Deep reinforcement learning with double Q-learning. In *AAAI conference on artificial intelligence*, 2016.

Appendix

A. κ -PI-DQN and κ -VI-DQN Algorithms

A.1. Detailed Pseudo-codes

In this section, we report the detailed pseudo-codes of κ -PI-DQN and κ -VI-DQN algorithms, described in Section 4.3, side-by-side.

Algorithm 5 κ -PI-DQN

```

1: Initialize replay buffer  $\mathcal{D}$ ;  $Q$ -networks  $Q_\theta$  and  $Q_\phi$  with random weights  $\theta$  and  $\phi$ ;
2: Initialize target networks  $Q'_\theta$  and  $Q'_\phi$  with weights  $\theta' \leftarrow \theta$  and  $\phi' \leftarrow \phi$ ;
3: for  $i = 0, \dots, N_\kappa - 1$  do
4:   # Policy Improvement
5:   for  $t = 1, \dots, T_\kappa$  do
6:     Select  $a_t$  as an  $\epsilon$ -greedy action w.r.t.  $Q_\theta(s_t, a)$ ;
7:     Execute  $a_t$ , observe  $r_t$  and  $s_{t+1}$ , and store the tuple  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ ;
8:     Sample a random mini-batch  $\{(s_j, a_j, r_j, s_{j+1})\}_{j=1}^N$  from  $\mathcal{D}$ ;
9:     Update  $\theta$  by minimizing the following loss function:
10:     $\mathcal{L}_{Q_\theta} = \frac{1}{N} \sum_{j=1}^N [Q_\theta(s_j, a_j) - (r_j(\kappa, V_\phi) + \gamma \kappa \max_a Q'_\theta(s_{j+1}, a))]^2$ , where
11:     $V_\phi(s_{j+1}) = Q_\phi(s_{j+1}, \pi_{i-1}(s_{j+1}))$  and  $\pi_{i-1}(s_{j+1}) \in \arg \max_a Q'_\theta(s_{j+1}, a)$ ;
12:    Copy  $\theta$  to  $\theta'$  occasionally ( $\theta' \leftarrow \theta$ );
13:   end for
14:   # Policy Evaluation
15:   Set  $\pi_i(s) \in \arg \max_a Q'_\theta(s, a)$ ;
16:   for  $t' = 1, \dots, T(\kappa)$  do
17:     Sample a random mini-batch  $\{(s_j, a_j, r_j, s_{j+1})\}_{j=1}^N$  from  $\mathcal{D}$ ;
18:     Update  $\phi$  by minimizing the following loss function:
19:     $\mathcal{L}_{Q_\phi} = \frac{1}{N} \sum_{j=1}^N [Q_\phi(s_j, a_j) - (r_j + \gamma Q'_\phi(s_{j+1}, \pi_i(s_{j+1})))]^2$ ;
20:    Copy  $\phi$  to  $\phi'$  occasionally ( $\phi' \leftarrow \phi$ );
21:   end for
22: end for

```

Algorithm 6 κ -VI-DQN

```

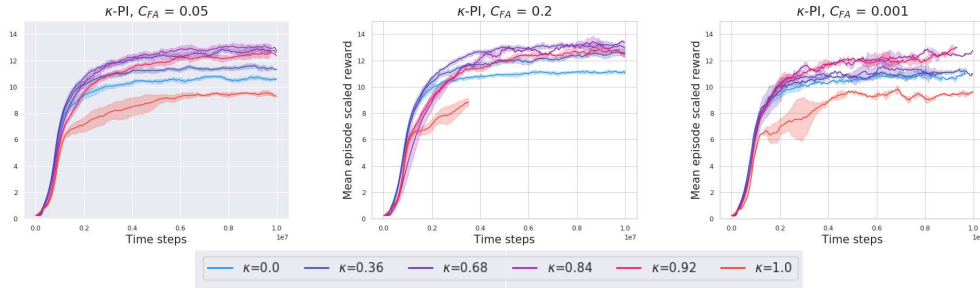
1: Initialize replay buffer  $\mathcal{D}$ ;  $Q$ -networks  $Q_\theta$  and  $Q_\phi$  with random weights  $\theta$  and  $\phi$ ;
2: Initialize target network  $Q'_\theta$  with weights  $\theta' \leftarrow \theta$ ;
3: for  $i = 0, \dots, N_\kappa - 1$  do
4:   # Evaluate  $T_\kappa V_\phi$  and the  $\kappa$ -greedy policy w.r.t.  $V_\phi$ 
5:   for  $t = 1, \dots, T_\kappa$  do
6:     Select  $a_t$  as an  $\epsilon$ -greedy action w.r.t.  $Q_\theta(s_t, a)$ ;
7:     Execute  $a_t$ , observe  $r_t$  and  $s_{t+1}$ , and store the tuple  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}$ ;
8:     Sample a random mini-batch  $\{(s_j, a_j, r_j, s_{j+1})\}_{j=1}^N$  from  $\mathcal{D}$ ;
9:     Update  $\theta$  by minimizing the following loss function:
10:     $\mathcal{L}_{Q_\theta} = \frac{1}{N} \sum_{j=1}^N [Q_\theta(s_j, a_j) - (r_j(\kappa, V_\phi) + \kappa \gamma \max_a Q'_\theta(s_{j+1}, a))]^2$ , where
11:     $V_\phi(s_{j+1}) = Q_\phi(s_{j+1}, \pi(s_{j+1}))$  and  $\pi(s_{j+1}) \in \arg \max_a Q_\phi(s_{j+1}, a)$ ;
12:    Copy  $\theta$  to  $\theta'$  occasionally ( $\theta' \leftarrow \theta$ );
13:   end for
14:   Copy  $\theta$  to  $\phi$  ( $\phi \leftarrow \theta$ )
15: end for

```

| Hyperparameter | Value |
|-----------------------------------|--------------------|
| Horizon (T) | 1000 |
| Adam stepsize | 1×10^{-4} |
| Target network update frequency | 1000 |
| Replay memory size | 100000 |
| Discount factor | 0.99 |
| Total training time steps | 10000000 |
| Minibatch size | 32 |
| Initial exploration | 1 |
| Final exploration | 0.1 |
| Final exploration frame | 1000000 |
| #Runs used for plot averages | 10 |
| Confidence interval for plot runs | $\sim 95\%$ |

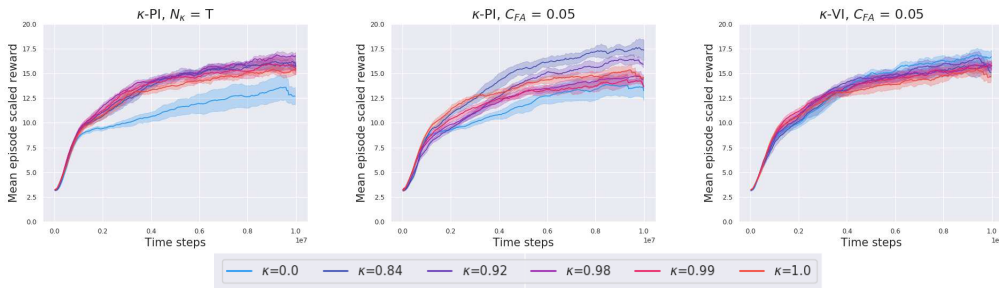
 Table 3: Hyperparameters for κ -PI-DQN and κ -VI-DQN.

A.2. Ablation Test for C_{FA}


 Figure 4: Performance of κ -PI-DQN and κ -VI-DQN on Breakout for different values of C_{FA} .

A.3. κ -PI-DQN and κ -VI-DQN Plots

In this section, we report additional results of the application of κ -PI-DQN and κ -VI-DQN on the Atari domains. A summary of these results has been reported in Table 1 in the main paper.


 Figure 5: Training performance of the 'naive' baseline $N_\kappa = T$ and κ -PI-DQN, κ -VI-DQN for $C_{FA} = 0.05$ on SpaceInvaders

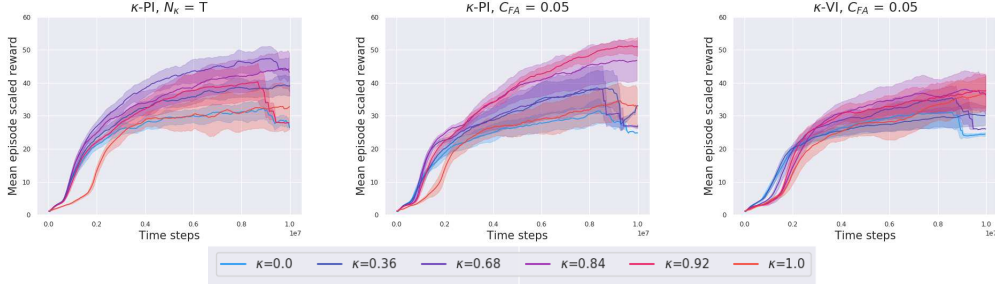


Figure 6: Training performance of the ‘naive’ baseline $N_\kappa = T$ and κ -PI-DQN, κ -VI-DQN for $C_{FA} = 0.05$ on Seaquest

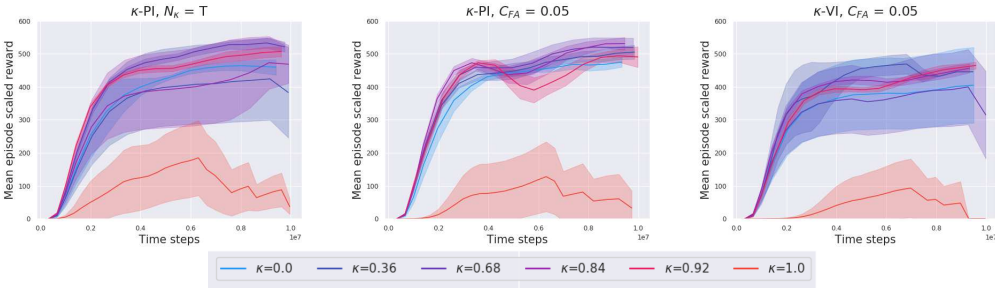


Figure 7: Training performance of the ‘naive’ baseline $N_\kappa = T$ and κ -PI-DQN, κ -VI-DQN for $C_{FA} = 0.05$ on Enduro

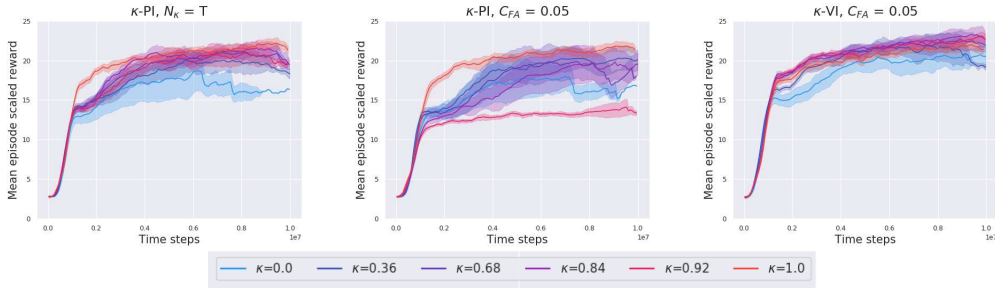


Figure 8: Training performance of the ‘naive’ baseline $N_\kappa = T$ and κ -PI-DQN, κ -VI-DQN for $C_{FA} = 0.05$ on BeamRider

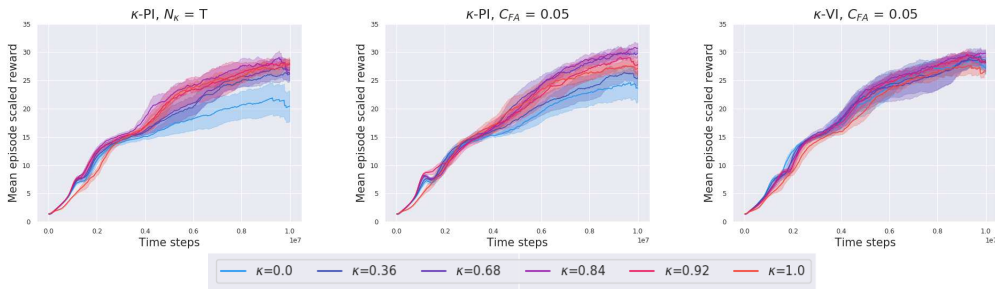


Figure 9: Training performance of the ‘naive’ baseline $N_\kappa = T$ and κ -PI-DQN, κ -VI-DQN for $C_{FA} = 0.05$ on Qbert

B. κ -PI-TRPO and κ -VI-TRPO Algorithms

B.1. Detailed Pseudo-codes

In this section, we report the detailed pseudo-codes of the κ -PI-TRPO and κ -VI-TRPO algorithms, described in Section 4.4, side-by-side.

Algorithm 7 κ -PI-TRPO

```

1: Initialize  $V$ -networks  $V_\theta$  and  $V_\phi$  with random weights  $\theta$  and  $\phi$ ; policy network  $\pi_\psi$  with random weights  $\psi$ ;
2: for  $i = 0, \dots, N_\kappa - 1$  do
3:   for  $t = 1, \dots, T_\kappa$  do
4:     Simulate the current policy  $\pi_\psi$  for  $M$  time-steps;
5:     for  $j = 1, \dots, M$  do
6:       Calculate  $R_j(\kappa, V_\phi) = \sum_{t=j}^M (\gamma\kappa)^{t-j} r_t(\kappa, V_\phi)$  and  $\rho_j = \sum_{t=j}^M \gamma^{t-j} r_t$ ;
7:     end for
8:     Sample a random mini-batch  $\{(s_j, a_j, r_j, s_{j+1})\}_{j=1}^N$  from the simulated  $M$  time-steps;
9:     Update  $\theta$  by minimizing the loss function:  $\mathcal{L}_{V_\theta} = \frac{1}{N} \sum_{j=1}^N (V_\theta(s_j) - R_j(\kappa, V_\phi))^2$ ;
10:    # Policy Improvement
11:    Sample a random mini-batch  $\{(s_j, a_j, r_j, s_{j+1})\}_{j=1}^N$  from the simulated  $M$  time-steps;
12:    Update  $\psi$  using TRPO with advantage function computed by  $\{(R_j(\kappa, V_\phi), V_\theta(s_j))\}_{j=1}^N$ ;
13:  end for
14:  # Policy Evaluation
15:  Sample a random mini-batch  $\{(s_j, a_j, r_j, s_{j+1})\}_{j=1}^N$  from the simulated  $M$  time-steps;
16:  Update  $\phi$  by minimizing the loss function:  $\mathcal{L}_{V_\phi} = \frac{1}{N} \sum_{j=1}^N (V_\phi(s_j) - \rho_j)^2$ ;
17: end for

```

Algorithm 8 κ -VI-TRPO

```

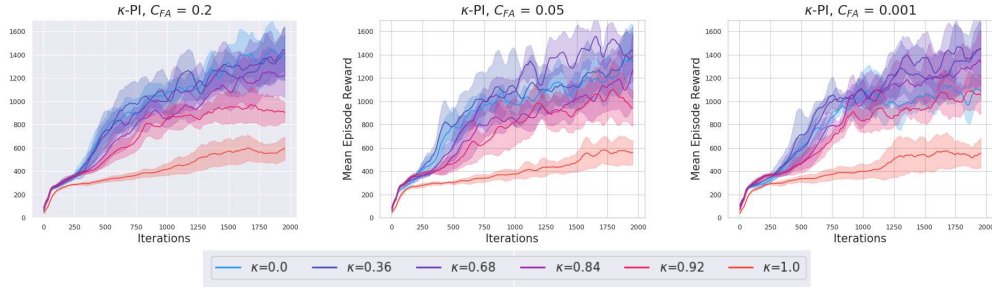
1: Initialize  $V$ -networks  $V_\theta$  and  $V_\phi$  with random weights  $\theta$  and  $\phi$ ; policy network  $\pi_\psi$  with random weights  $\psi$ ;
2: for  $i = 0, \dots, N_\kappa - 1$  do
3:   # Evaluate  $T_\kappa V_\phi$  and the  $\kappa$ -greedy policy w.r.t.  $V_\phi$ 
4:   for  $t = 1, \dots, T_\kappa$  do
5:     Simulate the current policy  $\pi_\psi$  for  $M$  time-steps;
6:     for  $j = 1, \dots, M$  do
7:       Calculate  $R_j(\kappa, V_\phi) = \sum_{t=j}^M (\gamma\kappa)^{t-j} r_t(\kappa, V_\phi)$ ;
8:     end for
9:     Sample a random mini-batch  $\{(s_j, a_j, r_j, s_{j+1})\}_{j=1}^N$  from the simulated  $M$  time-steps;
10:    Update  $\theta$  by minimizing the loss function:  $\mathcal{L}_{V_\theta} = \frac{1}{N} \sum_{j=1}^N (V_\theta(s_j) - R_j(\kappa, V_\phi))^2$ ;
11:    Sample a random mini-batch  $\{(s_j, a_j, r_j, s_{j+1})\}_{j=1}^N$  from the simulated  $M$  time-steps;
12:    Update  $\psi$  using TRPO with advantage function computed by  $\{(R_j(\kappa, V_\phi), V_\theta(s_j))\}_{j=1}^N$ ;
13:  end for
14:  Copy  $\theta$  to  $\phi$  ( $\phi \leftarrow \theta$ );
15: end for

```

| Hyperparameter | Value |
|-----------------------------------|--------------------|
| Horizon (T) | 1000 |
| Adam stepsize | 1×10^{-3} |
| Number of samples per Iteration | 1024 |
| Entropy coefficient | 0.01 |
| Discount factor | 0.99 |
| Number of Iterations | 2000 |
| Minibatch size | 128 |
| #Runs used for plot averages | 10 |
| Confidence interval for plot runs | $\sim 95\%$ |

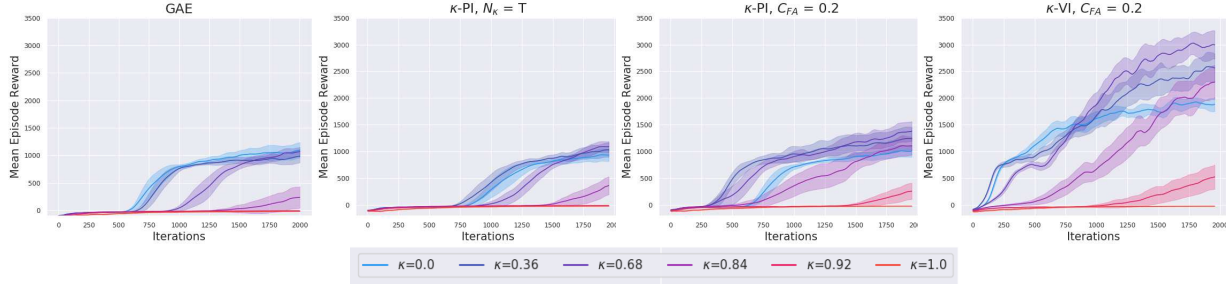
 Table 4: Hyper-parameters of κ -PI-TRPO and κ -VI-TRPO on the MuJoCo domains.

B.2. Ablation Test for C_{FA}


 Figure 10: Performance of κ -PI-TRPO and κ -VI-TRPO on Walker2d-v2 for different values of C_{FA} .

B.3. κ -PI-TRPO and κ -VI-TRPO Plots

In this section, we report additional results of the application of κ -PI-TRPO and κ -VI-TRPO on the MuJoCo domains. A summary of these results has been reported in Table 2 in the main paper.


 Figure 11: Performance of GAE, 'Naive' baseline and κ -PI-TRPO, κ -VI-TRPO on Ant-v2.

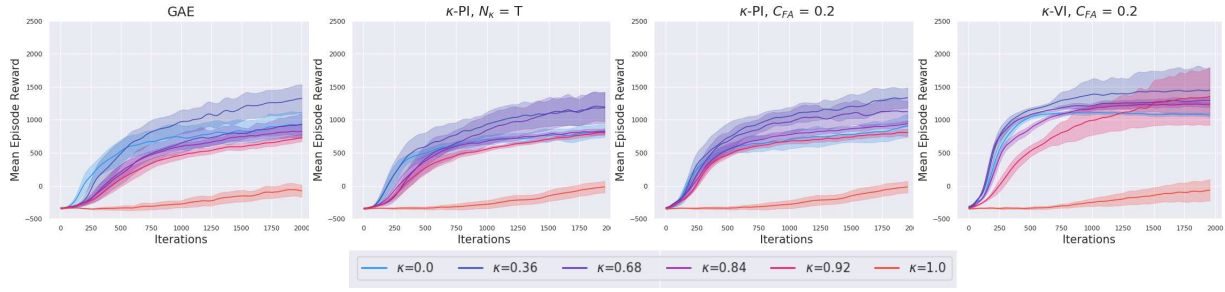


Figure 12: Performance of GAE, ‘Naive’ baseline and κ -PI-TRPO, κ -VI-TRPO on HalfCheetah-v2.

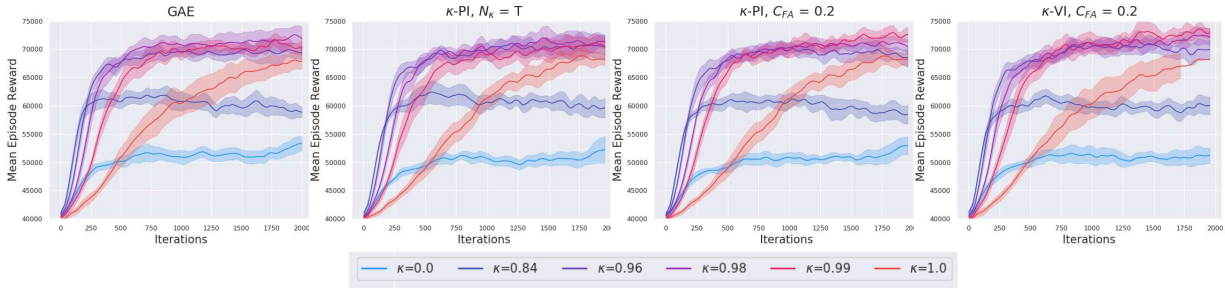


Figure 13: Performance of GAE, ‘Naive’ baseline and κ -PI-TRPO, κ -VI-TRPO on HumanoidStandup-v2.

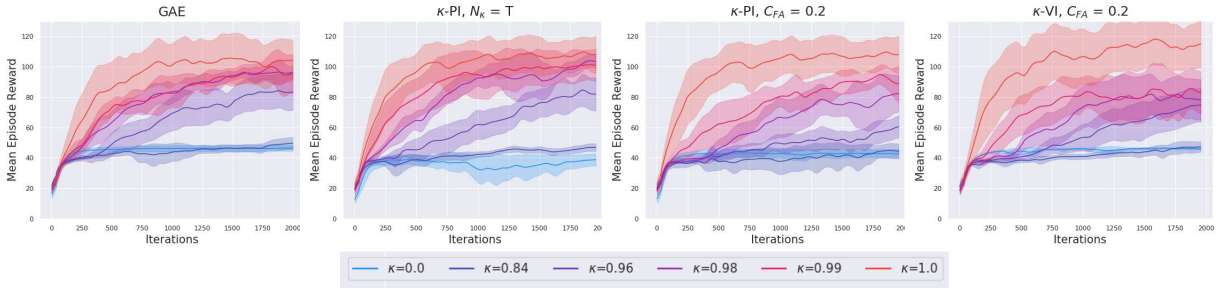


Figure 14: Performance of GAE, ‘Naive’ baseline and κ -PI-TRPO, κ -VI-TRPO on Swimmer-v2.

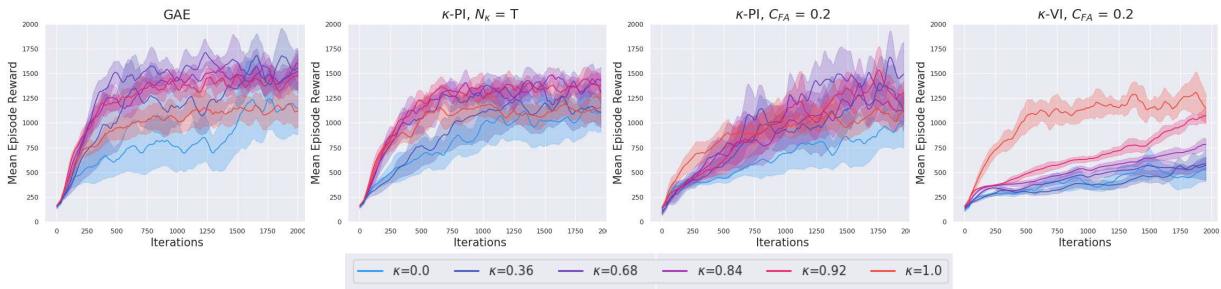


Figure 15: Performance of GAE, ‘Naive’ baseline and κ -PI-TRPO, κ -VI-TRPO on Hopper-v2.