

---

# On-Policy Deep Reinforcement Learning for the Average-Reward Criterion

---

Yiming Zhang<sup>1</sup> Keith W. Ross<sup>2,1</sup>

## Abstract

We develop theory and algorithms for average-reward on-policy Reinforcement Learning (RL). We first consider bounding the difference of the long-term average reward for two policies. We show that previous work based on the discounted return (Schulman et al., 2015; Achiam et al., 2017) results in a non-meaningful bound in the average-reward setting. By addressing the average-reward criterion directly, we then derive a novel bound which depends on the average divergence between the two policies and Kemeny’s constant. Based on this bound, we develop an iterative procedure which produces a sequence of monotonically improved policies for the average reward criterion. This iterative procedure can then be combined with classic DRL (Deep Reinforcement Learning) methods, resulting in practical DRL algorithms that target the long-run average reward criterion. In particular, we demonstrate that Average-Reward TRPO (ATRPO), which adapts the on-policy TRPO algorithm to the average-reward criterion, significantly outperforms TRPO in the most challenging MuJoCo environments.

## 1. Introduction

The goal of Reinforcement Learning (RL) is to build agents that can learn high-performing behaviors through trial-and-error interactions with the environment. Broadly speaking, modern RL tackles two kinds of problems: *episodic tasks* and *continuing tasks*. In episodic tasks, the agent-environment interaction can be broken into separate distinct episodes, and the performance of the agent is simply the sum of the rewards accrued within an episode. Examples of episodic tasks include training an agent to learn to play Go (Silver et al., 2016; 2018), where the episode terminates when the game ends. In continuing tasks, such as robotic locomotion (Peters & Schaal, 2008; Schulman et al., 2015;

Haarnoja et al., 2018) or in a queuing scenario (Tadepalli & Ok, 1994; Sutton & Barto, 2018), there is no natural separation of episodes and the agent-environment interaction continues indefinitely. The performance of an agent in a continuing task is more difficult to quantify since the total sum of rewards is typically infinite.

One way of making the long-term reward objective meaningful for continuing tasks is to apply *discounting* so that the infinite-horizon return is guaranteed to be finite for any bounded reward function. However the discounted objective biases the optimal policy to choose actions that lead to high near-term performance rather than to high long-term performance. Such an objective is not appropriate when the goal is to optimize long-term behavior, i.e., when the natural objective underlying the task at hand is non-discounted. In particular, we note that for the vast majority of benchmarks for reinforcement learning such as Atari games (Mnih et al., 2013) and MuJoCo (Todorov et al., 2012), a non-discounted performance measure is used to evaluate the trained policies.

Although in many circumstances, non-discounted criteria are more natural, most of the successful DRL algorithms today have been designed to optimize a discounted criterion during training. One possible work-around for this mismatch is to simply train with a discount factor that is very close to one. Indeed, from the Blackwell optimality theory of MDPs (Blackwell, 1962), we know that if the discount factor is very close to one, then an optimal policy for the infinite-horizon discounted criterion is also optimal for the long-run average-reward criterion. However, although Blackwell’s result suggests we can simply use a large discount factor to optimize non-discounted criteria, problems with large discount factors are in general more difficult to solve (Petrik & Scherrer, 2008; Jiang et al., 2015; 2016; Lehnert et al., 2018). Researchers have also observed that state-of-the-art DRL algorithms typically break down when the discount factor gets too close to one (Schulman et al., 2016; Andrychowicz et al., 2020).

In this paper we seek to develop algorithms for finding high-performing policies for average-reward DRL problems. Instead of trying to simply use standard discounted DRL algorithms with large discount factors, we instead attack the problem head-on, seeking to directly optimize the average-reward criterion. While the average reward setting has been

---

<sup>1</sup>New York University <sup>2</sup>New York University Shanghai. Correspondence to: Yiming Zhang <yiming.zhang@cs.nyu.edu>.

extensively studied in the classical Markov Decision Process literature (Howard, 1960; Blackwell, 1962; Veinott, 1966; Bertsekas et al., 1995), and has to some extent been studied for tabular RL (Schwartz, 1993; Mahadevan, 1996; Abounadi et al., 2001; Wan et al., 2020), it has received relatively little attention in the DRL community. In this paper, our focus is on developing average-reward on-policy DRL algorithms.

One major source of difficulty with modern on-policy DRL algorithms lies in controlling the step-size for policy updates. In order to have better control over step-sizes, Schulman et al. (2015) constructed a lower bound on the difference between the expected discounted return for two arbitrary policies  $\pi$  and  $\pi'$  by building upon the work of Kakade & Langford (2002). The bound is a function of the divergence between these two policies and the discount factor. Schulman et al. (2015) showed that iteratively maximizing this lower bound generates a sequence of monotonically improved policies for their discounted return.

In this paper, we first show that the policy improvement theorem from Schulman et al. (2015) results in a non-meaningful bound in the average reward case. We then derive a novel result which lower bounds the difference of the average long-run rewards. The bound depends on the average divergence between the policies and on the so-called Kemeny constant, which measures to what degree the irreducible Markov chains associated with the policies are “well-mixed”. We show that iteratively maximizing this lower bound guarantees monotonic average reward policy improvement.

Similar to the discounted case, the problem of maximizing the lower bound can be approximated with DRL algorithms which can be optimized using samples collected in the environment. In particular, we describe in detail the Average Reward TRPO (ATRPO) algorithm, which is the average reward variant of the TRPO algorithm (Schulman et al., 2015). Using the MuJoCo simulated robotic benchmark, we carry out extensive experiments demonstrating the effectiveness of ATRPO compared to its discounted counterpart, in particular on the most challenging MuJoCo tasks. Notably, we show that ATRPO can significantly out-perform TRPO on a set of high-dimensional continuing control tasks.

Our main contributions can be summarized as follows:

- We extend the policy improvement bound from Schulman et al. (2015) and Achiam et al. (2017) to the average reward setting. We demonstrate that our new bound depends on the average divergence between the two policies and on the mixing time of the underlying Markov chain.
- We use the aforementioned policy improvement bound to derive novel on-policy deep reinforcement learning algorithms for optimizing the average reward.
- Most modern DRL algorithms introduce a discount factor during training even when the natural objective of interest is undiscounted. This leads to a discrepancy between the evaluation and training objective. We demonstrate that optimizing the average reward directly can effectively address this mismatch and lead to much stronger performance.

## 2. Preliminaries

Consider a Markov Decision Process (MDP) (Sutton & Barto, 2018)  $(\mathcal{S}, \mathcal{A}, P, r, \mu)$  where the state space  $\mathcal{S}$  and action space  $\mathcal{A}$  are assumed to be finite. The transition probability is denoted by  $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ , the bounded reward function  $r : \mathcal{S} \times \mathcal{A} \rightarrow [r_{\min}, r_{\max}]$ , and  $\mu : \mathcal{S} \rightarrow [0, 1]$  is the initial state distribution. Let  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$  be a stationary policy where  $\Delta(\mathcal{A})$  is the probability simplex over  $\mathcal{A}$ , and  $\Pi$  is the set of all stationary policies. We consider two classes of MDPs:

**Assumption 1 (Ergodic).** *For every stationary policy, the induced Markov chain is irreducible and aperiodic.*

**Assumption 2 (Aperiodic Unichain).** *For every stationary policy, the induced Markov chain contains a single aperiodic recurrent class and a finite but possibly empty set of transient states.*

By definition, any MDP which satisfies Assumption 1 is also unichain. We note that most MDPs of practical interest belong in these two classes. We will mostly focus on MDPs which satisfy Assumption 1 in the main text. In the supplementary material, we will address the aperiodic unichain case. Here we present the two objective formulations for continuing control tasks: the average reward approach and discounted reward criterion.

### Average Reward Criterion

The average reward objective is defined as:

$$\rho(\pi) := \lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{N-1} r(s_t, a_t) \right] = \mathbb{E}_{\substack{s \sim d_\pi \\ a \sim \pi}} [r(s, a)]. \quad (1)$$

Here  $d_\pi(s) := \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=0}^{N-1} P_{\tau \sim \pi}(s_t = s)$  is the stationary state distribution under policy  $\pi$ , and  $\tau = (s_0, a_0, \dots)$  is a sample trajectory. The limits in  $\rho(\pi)$  and  $d_\pi(s)$  are guaranteed to exist under our assumptions. Since the MDP is aperiodic, it can also be shown that  $d_\pi(s) = \lim_{t \rightarrow \infty} P_{\tau \sim \pi}(s_t = s)$ . In the unichain case, the average reward  $\rho(\pi)$  does not depend on the initial state for any policy  $\pi$  (Bertsekas et al., 1995). We express the average-reward bias function as

$$\bar{V}^\pi(s) := \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} (r(s_t, a_t) - \rho(\pi)) \middle| s_0 = s \right]$$

and *average-reward action-bias function* as

$$\bar{Q}^\pi(s, a) := \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} (r(s_t, a_t) - \rho(\pi)) \middle| s_0 = s, a_0 = a \right].$$

We define the *average-reward advantage function* as

$$\bar{A}^\pi(s, a) := \bar{Q}^\pi(s, a) - \bar{V}^\pi(s).$$

### Discounted Reward Criterion

For some discount factor  $\gamma \in (0, 1)$ , the discounted reward objective is defined as

$$\rho_\gamma(\pi) := \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] = \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d_{\pi, \gamma} \\ a \sim \pi}} [r(s, a)] \quad (2)$$

where  $d_{\pi, \gamma}(s) := (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P_{\tau \sim \pi}(s_t = s)$  is known as the *future discounted state visitation distribution under policy*  $\pi$ . Note that unlike the average reward objective, the discounted objective depends on the initial state distribution  $\mu$ . It can be easily shown that  $d_{\pi, \gamma}(s) \rightarrow d_\pi(s)$  for all  $s$  as  $\gamma \rightarrow 1$ . The *discounted value function* is defined as  $V_\gamma^\pi(s) := \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0 = s \right]$  and *discounted action-value function*  $Q_\gamma^\pi(s, a) := \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0 = s, a_0 = a \right]$ . Finally, the *discounted advantage function* is defined as  $A_\gamma^\pi(s, a) := Q_\gamma^\pi(s, a) - V_\gamma^\pi(s)$ .

It is well-known that  $\lim_{\gamma \rightarrow 1} (1 - \gamma) \rho_\gamma(\pi) = \rho(\pi)$ , implying that the discounted and average reward objectives are equivalent in the limit as  $\gamma$  approaches 1 (Blackwell, 1962). We further discuss the relationship between the discounted and average reward criteria in Appendix A and prove that  $\lim_{\gamma \rightarrow 1} A_\gamma^\pi(s, a) = \bar{A}^\pi(s, a)$  (see Corollary A.1). The proofs of all results in the subsequent sections, if not given, can be found in the supplementary material.

### 3. Monotonically Improvement Guarantees for Discounted RL

In much of the on-policy DRL literature (Schulman et al., 2015; 2017; Wu et al., 2017; Vuong et al., 2019; Song et al., 2020), algorithms iteratively update policies by maximizing them within a local region, i.e., at iteration  $k$  we find a policy  $\pi_{k+1}$  by maximizing  $\rho_\gamma(\pi)$  within some region  $D(\pi, \pi_k) \leq \delta$  for some divergence measure  $D$ . By using different choices of  $D$  and  $\delta$ , this approach allows us to control the step-size of each update, which can lead to better sample efficiency (Peters & Schaal, 2008). Schulman et al. (2015) derived a policy improvement bound based on

a specific choice of  $D$ :

$$\begin{aligned} \rho_\gamma(\pi_{k+1}) - \rho_\gamma(\pi_k) &\geq \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d_{\pi_k, \gamma} \\ a \sim \pi_{k+1}}} [A_\gamma^{\pi_k}(s, a)] \\ &\quad - C \cdot \max_s [D_{\text{TV}}(\pi_{k+1} \parallel \pi_k)[s]] \end{aligned} \quad (3)$$

where  $D_{\text{TV}}(\pi' \parallel \pi)[s] := \frac{1}{2} \sum_a |\pi'(a|s) - \pi(a|s)|$  is the *total variation divergence*, and  $C = 4\gamma\epsilon/(1 - \gamma)^2$  where  $\epsilon$  is some constant. Schulman et al. (2015) showed that by choosing  $\pi_{k+1}$  which maximizes the right hand side of (3), we are guaranteed to have  $\rho_\gamma(\pi_{k+1}) \geq \rho_\gamma(\pi_k)$ . This provided the theoretical foundation for an entire class of on-policy DRL algorithms (Schulman et al., 2015; 2017; Wu et al., 2017; Vuong et al., 2019; Song et al., 2020).

A natural question arises here is whether the iterative procedure described by Schulman et al. (2015) also guarantees improvement for the average reward. Since the discounted and average reward objectives become equivalent as  $\gamma \rightarrow 1$ , one may conjecture that we can also lower bound the policy performance difference of the average reward objective by simply letting  $\gamma \rightarrow 1$  for the bounds in Schulman et al. (2015). Unfortunately this results in a non-meaningful bound (see supplementary material for proof.)

**Proposition 1.** *Consider the bounds in Theorem 1 of Schulman et al. (2015) and Corollary 1 of Achiam et al. (2017). The right hand side of both bounds times  $1 - \gamma$  goes to negative infinity as  $\gamma \rightarrow 1$ .*

Since  $\lim_{\gamma \rightarrow 1} (1 - \gamma)(\rho_\gamma(\pi') - \rho_\gamma(\pi)) = \rho(\pi') - \rho(\pi)$ , Proposition 1 says that the policy improvement guarantee from Schulman et al. (2015) and Achiam et al. (2017) becomes trivial when  $\gamma \rightarrow 1$  and thus does not generalize to the average reward setting. In the next section, we will derive a novel policy improvement bound for the average reward objective, which in turn can be used to generate monotonically improved policies w.r.t. the average reward.

## 4. Main Results

### 4.1. Average Reward Policy Improvement Theorem

Let  $d_\pi \in \mathbb{R}^{|S|}$  be the probability column vector whose components are  $d_\pi(s)$ . Let  $P_\pi \in \mathbb{R}^{|S| \times |S|}$  be the transition matrix under policy  $\pi$  whose  $(s, s')$  component is  $P_\pi(s'|s) = \sum_a P(s'|s, a)\pi(a|s)$ , and  $P_\pi^* := \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=0}^N P_\pi^t$  be the limiting distribution of the transition matrix. For aperiodic unichain MDPs,  $P_\pi^* = \lim_{t \rightarrow \infty} P_\pi^t = \mathbf{1}d_\pi^T$ .

Suppose we have a new policy  $\pi'$  obtained via some update rule from the current policy  $\pi$ . Similar to the discounted case, we would like to measure their performance difference  $\rho(\pi') - \rho(\pi)$  using an expression which depends on  $\pi$  and some divergence metric between the two policies. The following identity shows that  $\rho(\pi') - \rho(\pi)$  can be expressed

using the average reward advantage function of  $\pi$ .

**Lemma 1.** *Under Assumption 2:*

$$\rho(\pi') - \rho(\pi) = \mathbb{E}_{\substack{s \sim d_{\pi'} \\ a \sim \pi'}} [\bar{A}^\pi(s, a)] \quad (4)$$

for any two stochastic policies  $\pi$  and  $\pi'$ .

Lemma 1 is an extension of the well-known policy difference lemma from Kakade & Langford (2002) to the average reward case. A similar result was proven by Even-Dar et al. (2009) and Neu et al. (2010). For completeness, we provide a simple proof in the supplementary material. Note that this expression depends on samples drawn from  $\pi'$ . However we can show through the following lemma that when  $d_\pi$  and  $d_{\pi'}$  are ‘‘close’’ w.r.t. the TV divergence, we can evaluate  $\rho(\pi')$  using samples from  $d_\pi$  (see supplementary material for proof).

**Lemma 2.** *Under Assumption 2, the following bound holds for any two stochastic policies  $\pi$  and  $\pi'$ :*

$$\left| \rho(\pi') - \rho(\pi) - \mathbb{E}_{\substack{s \sim d_\pi \\ a \sim \pi'}} [\bar{A}^\pi(s, a)] \right| \leq 2\epsilon D_{TV}(d_{\pi'} \parallel d_\pi) \quad (5)$$

where  $\epsilon = \max_s |\mathbb{E}_{a \sim \pi'(a|s)} [\bar{A}^\pi(s, a)]|$ .

Lemma 2 implies that

$$\rho(\pi') \approx \rho(\pi) + \mathbb{E}_{\substack{s \sim d_\pi \\ a \sim \pi'}} [\bar{A}^\pi(s, a)] \quad (6)$$

when  $d_\pi$  and  $d_{\pi'}$  are ‘‘close’’. However in order to study how policy improvement is connected to changes in the actual policies themselves, we need to analyze the relationship between changes in the policies and changes in stationary distributions. It turns out that the sensitivity of the stationary distributions in relation to the policies is related to the structure of the underlying Markov chain.

Let  $M^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$  be the *mean first passage time matrix* whose elements  $M^\pi(s, s')$  is the expected number of steps it takes to reach state  $s'$  from  $s$  under policy  $\pi$ . Under Assumption 1, the matrix  $M^\pi$  can be calculated via (see Theorem 4.4.7 of Kemeny & Snell (1960))

$$M^\pi = (I - Z^\pi + EZ_{\text{dg}}^\pi)D^\pi \quad (7)$$

where  $Z^\pi = (I - P_\pi + P_\pi^*)^{-1}$  is known as the *fundamental matrix of the Markov chain* (Kemeny & Snell, 1960),  $E$  is a square matrix consisting of all ones. The subscript ‘‘dg’’ on some square matrix refers to taking the diagonal of said matrix and placing zeros everywhere else.  $D^\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$  is a diagonal matrix whose elements are  $1/d_\pi(s)$ .

One important property of mean first passage time is that for any MDP which satisfies Assumption 1, the quantity

$$\kappa^\pi = \sum_{s'} d_\pi(s') M^\pi(s, s') = \text{trace}(Z^\pi) \quad (8)$$

is a constant independent of the starting state for any policy  $\pi$  (Theorem 4.4.10 of Kemeny & Snell (1960).) The constant  $\kappa^\pi$  is sometimes referred to as *Kemeny’s constant* (Grinstead & Snell, 2012). This constant can be interpreted as the mean number of steps it takes to get to any goal state weighted by the steady-distribution of the goal states. This weighted mean does not depend on the starting state, as mentioned just above.

It can be shown that the value of Kemeny’s constant is also related to the *mixing time* of the Markov Chain, i.e., how fast the chain converges to the stationary distribution (see Appendix C for additional details).

The following result connects the sensitivity of the stationary distribution to changes to the policy.

**Lemma 3.** *Under Assumption 1, the divergence between the stationary distributions  $d_\pi$  and  $d_{\pi'}$  can be upper bounded by the average divergence between policies  $\pi$  and  $\pi'$ :*

$$D_{TV}(d_{\pi'} \parallel d_\pi) \leq (\kappa^* - 1) \mathbb{E}_{s \sim d_\pi} [D_{TV}(\pi' \parallel \pi)[s]] \quad (9)$$

where  $\kappa^* = \max_\pi \kappa^\pi$

For Markov chains with a small mixing time, where an agent can quickly get to any state, Kemeny’s constant is relatively small and Lemma 3 shows that the stationary distributions are not highly sensitive to small changes in the policy. On the other hand, for Markov chains that have high mixing times, the factor can become very large. In this case Lemma 3 shows that small changes in the policy can have a large impact on the resulting stationary distributions.

Combining the bounds in Lemma 2 and Lemma 3 gives us the following result:

**Theorem 1.** *Under Assumption 1 the following bounds hold for any two stochastic policies  $\pi$  and  $\pi'$ , :*

$$D_\pi^-(\pi') \leq \rho(\pi') - \rho(\pi) \leq D_\pi^+(\pi') \quad (10)$$

where

$$D_\pi^\pm(\pi') = \mathbb{E}_{\substack{s \sim d_\pi \\ a \sim \pi'}} [\bar{A}^\pi(s, a)] \pm 2\xi \mathbb{E}_{s \sim d_\pi} [D_{TV}(\pi' \parallel \pi)[s]]$$

and  $\xi = (\kappa^* - 1) \max_s \mathbb{E}_{a \sim \pi'} |\bar{A}^\pi(s, a)|$ .

The bounds in Theorem 1 are guaranteed to be finite. Analogous to the discounted case, the multiplicative factor  $\xi$  provides guidance on the step-sizes for policy updates. Note that Theorem 1 holds for MDPs satisfying Assumption 1; in Appendix D we discuss how a similar result can be derived for the more general aperiodic unichain case.

The bound in Theorem 1 is given in terms of the TV divergence; however the KL divergence is more commonly used in practice. The relationship between the TV divergence and



**Algorithm 1** Approximate Average Reward Policy Iteration

- 1: **Input:**  $\pi_0$
- 2: **for**  $k = 0, 1, 2, \dots$  **do**
- 3:   Policy Evaluation: Evaluate  $\bar{A}^{\pi_k}(s, a)$  for all  $s, a$
- 4:   Policy Improvement:

$$\pi_{k+1} = \operatorname{argmax}_{\pi} D_{\pi_k}^-(\pi) \quad (12)$$

where

$$D_{\pi_k}^-(\pi) = \mathbb{E}_{\substack{s \sim d_{\pi_k} \\ a \sim \pi}} [\bar{A}^{\pi_k}(s, a)] - \xi \sqrt{2 \mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi \| \pi_k)[s]]}$$

and  $\xi = (\kappa^* - 1) \max_s \mathbb{E}_{a \sim \pi} |\bar{A}^{\pi_k}(s, a)|$

5: **end for**

KL divergence is given by Pinsker's inequality (Tsybakov, 2008), which says that for any two distributions  $p$  and  $q$ :  $D_{\text{TV}}(p \| q) \leq \sqrt{D_{\text{KL}}(p \| q)}/2$ . We can then show that

$$\begin{aligned} \mathbb{E}_{s \sim d_{\pi}} [D_{\text{TV}}(\pi' \| \pi)[s]] &\leq \mathbb{E}_{s \sim d_{\pi}} [\sqrt{D_{\text{KL}}(\pi' \| \pi)[s]}/2] \\ &\leq \sqrt{\mathbb{E}_{s \sim d_{\pi}} [D_{\text{KL}}(\pi' \| \pi)[s] ]}/2 \end{aligned} \quad (11)$$

where the second inequality comes from Jensen's inequality. The inequality in (11) shows that the bounds in Theorem 1 still hold when  $\mathbb{E}_{s \sim d_{\pi}} [D_{\text{TV}}(\pi' \| \pi)[s]]$  is substituted with  $\sqrt{\mathbb{E}_{s \sim d_{\pi}} [D_{\text{KL}}(\pi' \| \pi)[s] ]}/2$ .

## 4.2. Approximate Policy Iteration

One direct consequence of Theorem 1 is that iteratively maximizing the  $D_{\pi}^-(\pi')$  term in the bound generates a monotonically improving sequence of policies w.r.t. the average reward objective. Algorithm 1 gives an approximate policy iteration algorithm that produces such a sequence of policies.

**Proposition 2.** *Given an initial policy  $\pi_0$ , Algorithm 1 is guaranteed to generate a sequence of policies  $\pi_1, \pi_2, \dots$  such that  $\rho(\pi_0) \leq \rho(\pi_1) \leq \rho(\pi_2) \leq \dots$ .*

*Proof.* At iteration  $k$ ,  $\mathbb{E}_{s \sim d_{\pi_k}, a \sim \pi} [\bar{A}^{\pi_k}(s, a)] = 0$ ,  $\mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi \| \pi_k)[s]] = 0$  for  $\pi = \pi_k$ . By Theorem 1 and (12),  $\rho(\pi_{k+1}) - \rho(\pi_k) \geq 0$ .  $\square$

However, Algorithm 1 is difficult to implement in practice since it requires exact knowledge of  $\bar{A}^{\pi_k}(s, a)$  and the transition matrix. Furthermore, calculating the term  $\xi$  is impractical for high-dimensional problems. In the next section, we will introduce a sample-based algorithm which approximates the update rule in Algorithm 1.

## 5. Practical Algorithm

As noted in the previous section, Algorithm 1 is not practical for problems with large state and action spaces. In this section, we will discuss how Algorithm 1 and Theorem 1 can be used in practice to create algorithms which can effectively solve high dimensional DRL problems with the use of *trust region* methods.

In Appendix F, we will also discuss how Theorem 1 can be used to solve DRL problems with average cost safety constraints. RL with safety constraints are an important class of problems with practical implications (Amodei et al., 2016). Trust region methods have been successfully applied to this class of problems as it provides worst-case constraint violation guarantees for evaluating the cost constraint values for policy updates (Achiam et al., 2017; Yang et al., 2020; Zhang et al., 2020). However the aforementioned theoretical guarantees were only shown to apply to discounted cost constraints. Tessler et al. (2019) pointed out that trust-region based methods such as the Constrained Policy Optimization (CPO) algorithm (Achiam et al., 2017) cannot be used for average costs constraints. Contrary to this belief, in Appendix F, we demonstrate that Theorem 1 provides a worst-case constraint violation guarantee for average costs and trust-region-based constrained RL methods can easily be modified to accommodate for average cost constraints.

### 5.1. Average Reward Trust Region Methods

For DRL problems, it is common to consider some parameterized policy class  $\Pi_{\Theta} = \{\pi_{\theta} : \theta \in \Theta\}$ . Our goal is to devise a computationally tractable version of Algorithm 1 for policies in  $\Pi_{\Theta}$ . We can rewrite the unconstrained optimization problem in (12) as a constrained problem:

$$\begin{aligned} &\underset{\pi_{\theta} \in \Pi_{\Theta}}{\text{maximize}} && \mathbb{E}_{\substack{s \sim d_{\pi_{\theta_k}} \\ a \sim \pi_{\theta}}} [\bar{A}^{\pi_{\theta_k}}(s, a)] \\ &\text{subject to} && \bar{D}_{\text{KL}}(\pi_{\theta} \| \pi_{\theta_k}) \leq \delta \end{aligned} \quad (13)$$

where  $\bar{D}_{\text{KL}}(\pi_{\theta} \| \pi_{\theta_k}) := \mathbb{E}_{s \sim d_{\pi_{\theta_k}}} [D_{\text{KL}}(\pi_{\theta} \| \pi_{\theta_k})[s]]$ . Importantly, the advantage function  $\bar{A}^{\pi_{\theta_k}}(s, a)$  appearing in (13) is the average-reward advantage function, defined as the bias minus the action-bias, and not the discounted advantage function. The constraint set  $\{\pi_{\theta} \in \Pi_{\Theta} : \bar{D}_{\text{KL}}(\pi_{\theta} \| \pi_{\theta_k}) \leq \delta\}$  is called the *trust region set*. The problem (13) can be regarded as an average reward variant of the trust region problem from Schulman et al. (2015). The step-size  $\delta$  is treated as a hyperparameter in practice and should ideally be tuned for each specific task. However we note that in the average reward, the choice of step-size is related to the mixing time of the underlying Markov chain (since it is related to the multiplicative factor  $\xi$  in Theorem 1). When the mixing time is small, a larger step-size can be chosen and vice versa. While it is impractical to calculate the optimal

step-size, in certain applications domain knowledge on the mixing time can be used to serve as a guide for tuning  $\delta$ .

When we set  $\pi_{\theta_{k+1}}$  to be the optimal solution to (13), similar to the discounted case, the policy improvement guarantee no longer holds. However we can show that  $\pi_{\theta_{k+1}}$  has the following worst-case performance degradation guarantee:

**Proposition 3.** *Let  $\pi_{\theta_{k+1}}$  be the optimal solution to (13) for some  $\pi_{\theta_k} \in \Pi_{\Theta}$ . The policy performance difference between  $\pi_{\theta_{k+1}}$  and  $\pi_{\theta_k}$  can be lower bounded by*

$$\rho(\pi_{\theta_{k+1}}) - \rho(\pi_{\theta_k}) \geq -\xi^{\pi_{\theta_{k+1}}} \sqrt{2\delta} \quad (14)$$

where  $\xi^{\pi_{\theta_{k+1}}} = (\kappa^{\pi_{\theta_{k+1}}} - 1) \max_s \mathbb{E}_{a \sim \pi_{\theta_{k+1}}} |\bar{A}^{\pi_{\theta_k}}(s, a)|$ .

*Proof.* Since  $\bar{D}_{\text{KL}}(\pi_{\theta_k} \parallel \pi_{\theta_k}) = 0$ ,  $\pi_{\theta_k}$  is feasible. The objective value is 0 for  $\pi_{\theta} = \pi_{\theta_k}$ . The bound follows from (10) and (11) where the average KL is bounded by  $\delta$ .  $\square$

Several algorithms have been proposed for efficiently solving the discounted version of (13): Schulman et al. (2015) and Wu et al. (2017) converts (13) into a convex problem via Taylor approximations; another approach is to first solve (13) in the non-parametric policy space and then project the result back into the parameter space (Vuong et al., 2019; Song et al., 2020). These algorithms can also be adapted for the average reward case and are theoretically justified via Theorem 1 and Proposition 3. In the next section, we will provide as a specific example how this can be done for one such algorithm.

## 5.2. Average Reward TRPO (ATRPO)

In this section, we introduce ATRPO, which is an average-reward modification of the TRPO algorithm (Schulman et al., 2015). Similar to TRPO, we apply Taylor approximations to (13). This gives us a new optimization problem which can be solved exactly using Lagrange duality (Boyd et al., 2004). The solution to this approximate problem gives an explicit update rule for the policy parameters which then allows us to perform policy updates using an actor-critic framework. More details can be found in Appendix E. Algorithm 2 provides a basic outline of ATRPO.

The major differences between ATRPO and TRPO are as follows:

- i The critic network in Algorithm 2 approximates the average-reward bias rather than the discounted value function.
- ii ATRPO must estimate the average return  $\rho$  of the current policy.
- iii The targets for the bias and the advantage are calculated without discount factors and the average return  $\rho$  is

---

### Algorithm 2 Average Reward TRPO (ATRPO)

---

- 1: **Input:** Policy parameters  $\theta_0$ , critic net parameters  $\phi_0$ , learning rate  $\alpha$ , trajectory truncation parameter  $N$ .
- 2: **for**  $k = 0, 1, 2, \dots$  **do**
- 3:   Collect a truncated trajectory  $\{s_t, a_t, s_{t+1}, r_t\}$ ,  $t = 1, \dots, N$  from the environment using  $\pi_{\theta_k}$ .
- 4:   Calculate sample average reward of  $\pi_{\theta_k}$  via  $\rho = \frac{1}{N} \sum_{t=1}^N r_t$ .
- 5:   **for**  $t = 1, 2, \dots, N$  **do**
- 6:     Get target  $\bar{V}_t^{\text{target}} = r_t - \rho + \bar{V}_{\phi_k}(s_{t+1})$
- 7:     Get advantage estimate:  $\hat{A}(s_t, a_t) = r_t - \rho + \bar{V}_{\phi_k}(s_{t+1}) - \bar{V}_{\phi_k}(s_t)$
- 8:   **end for**
- 9:   Update critic by

$$\phi_{k+1} \leftarrow \phi_k - \alpha \nabla_{\phi} \mathcal{L}(\phi_k)$$

where

$$\mathcal{L}(\phi_k) = \frac{1}{N} \sum_{t=1}^N \|\bar{V}_{\phi_k}(s_t) - \bar{V}_t^{\text{target}}\|^2$$

- 10:   Use  $\hat{A}(s_t, a_t)$  to update  $\theta_k$  using TRPO policy update (Schulman et al., 2015).
  - 11: **end for**
- 

subtracted from the reward. Simply setting the discount factor to 1 in TRPO does not lead to Algorithm 2.

- iv ATRPO also assumes that the underlying task is a continuing infinite-horizon task. But since in practice we cannot run infinitely long trajectories, all trajectories are truncated at some large truncation value  $N$ . Unlike TRPO, during training we do not allow for episodic tasks where episodes terminate early (before  $N$ ). For the MuJoCo environments, we will address this by having the agent not only resume locomotion after falling but also incur a penalty for falling (see Section 6).

In Algorithm 2, for illustrative purposes, we use the average reward one-step bootstrapped estimate for the target of the critic and the advantage function. In practice, we instead develop and use an average-reward version of the Generalized Advantage Estimator (GAE) from Schulman et al. (2016). In Appendix G we provide more details on how GAE can be generalized to the average-reward case.

## 6. Experiments

We conducted experiments comparing the performance of ATRPO and TRPO on continuing control tasks. We consider three tasks (Ant, HalfCheetah, and Humanoid) from the MuJoCo physical simulator (Todorov et al., 2012) imple-

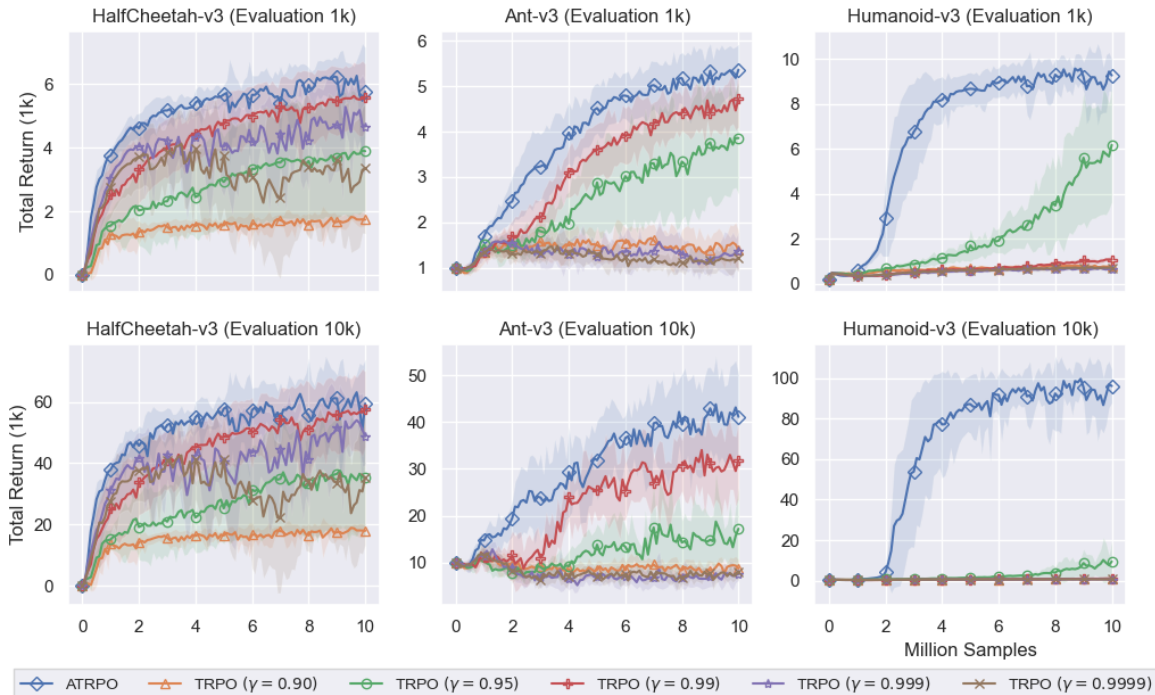


Figure 1. Comparing performance of ATRPO and TRPO with different discount factors. The  $x$ -axis is the number of agent-environment interactions and the  $y$ -axis is the total return averaged over 10 seeds. The solid line represents the agents’ performance on evaluation trajectories of maximum length 1,000 (top row) and 10,000 (bottom row). The shaded region represents one standard deviation.

mented using OpenAI gym (Brockman et al., 2016), where the natural goal is to train the agents to run as fast as possible without falling.

### 6.1. Evaluation Protocol

Even though the MuJoCo benchmark is commonly trained using the *discounted* objective (see e.g. Schulman et al. (2015), Wu et al. (2017), Lillicrap et al. (2016), Schulman et al. (2017), Haarnoja et al. (2018), Vuong et al. (2019)), it is *always* evaluated without discounting. Similarly, we also evaluate performance using the undiscounted total-reward objective for both TRPO and ATRPO.

Specifically for each environment, we train a policy for 10 million environment steps. During training, every 100,000 steps, we run 10 separate evaluation trajectories with the current policy without exploration (i.e., the policy is kept fixed and deterministic). For each evaluation trajectory we calculate the undiscounted return of the trajectory until the agent falls or until 1,000 steps, whichever comes first. We then report the average undiscounted return over the 10 trajectories. *Note that this is the standard evaluation metric for the MuJoCo environments.* In order to understand the performance of the agent for long time horizons, we also report the performance of the agent evaluated on trajectories of maximum length 10,000.

### 6.2. Comparing ATRPO and TRPO

To simulate an infinite-horizon setting during training, we do the following: when the agent falls, the trajectory does not terminate; instead the agent incurs a large reset cost for falling, and then continues the trajectory from a random start state. The reset cost is set to 100. However, we show in the supplementary material (Appendix I.2) that the results are largely insensitive to the choice of reset cost. We note that this modification does not change the underlying goal of the task. We also point out that the reset cost is only applied during training and is not used in the evaluation phase described in the previous section. Hyperparameter settings and other additional details can be found in Appendix H.

We plot the performance for ATRPO and TRPO trained with different discount factors in Figure 1. We see that TRPO with its best discount factor can perform as well as ATRPO for the simplest environment HalfCheetah. But ATRPO provides dramatic improvements in Ant and Humanoid. In particular for the most challenging environment Humanoid, ATRPO performs on average 50.1% better than TRPO with its best discount factor when evaluated on trajectories of maximum length 1000. The improvement is even greater when the agents are evaluated on trajectories of maximum length 10,000 where the performance boost jumps to 913%. In Appendix I.1, we provide an additional set of experiments

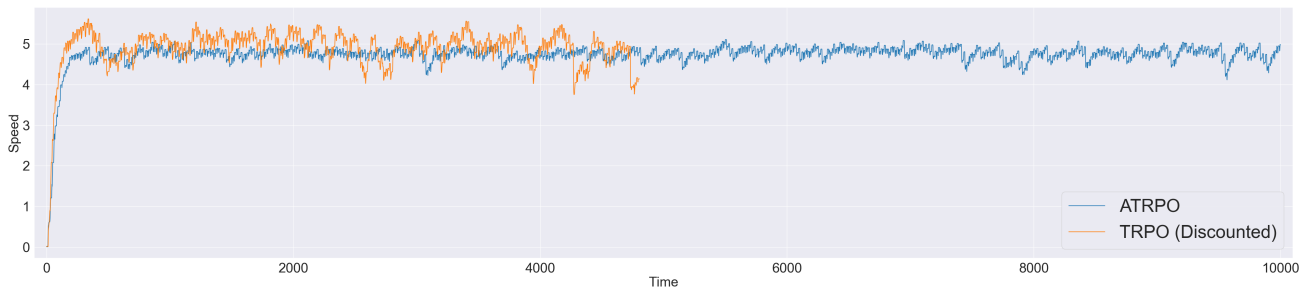


Figure 2. Speed-time plot of a single trajectory (maximum length 10,000) for ATRPO and Discounted TRPO in the Humanoid-v3 environment. The solid line represents the speed of the agent at the corresponding timesteps.

demonstrating that ATRPO also significantly outperforms TRPO when TRPO is trained without the reset scheme described at the beginning of this section (i.e. the standard MuJoCo setting.)

We make two observations regarding discounting. First, we note that increasing the discount factor does not necessarily lead to better performance for TRPO. A larger discount factor in principle enables the algorithm to seek a policy that performs well for the average-reward criterion (Blackwell, 1962). Unfortunately, a larger discount factor can also increase the variance of the gradient estimator (Zhao et al., 2011; Schulman et al., 2016), increase the complexity of the policy space (Jiang et al., 2015), lead to slower convergence (Bertsekas et al., 1995; Agarwal et al., 2020), and degrade generalization in limited data settings (Amit et al., 2020). Moreover, algorithms with discounting are known to become unstable as  $\gamma \rightarrow 1$  (Naik et al., 2019). Secondly, for TRPO the best discount factor is different for each environment (0.99 for HalfCheetah and Ant, 0.95 for Humanoid). The discount factor therefore serves as a hyperparameter which can be tuned to improve performance, choosing a suboptimal discount factor can have significant consequences. Both of these observations are consistent with what was seen in the literature (Andrychowicz et al., 2020). We have shown here that using the average reward criterion directly not only delivers superior performance but also obviates the need to tune the discount factor.

### 6.3. Understanding Long Run Performance

Next, we demonstrate that agents trained using the average reward criterion are better at optimizing for long-term returns. Here, we first train Humanoid with 10 million samples with ATRPO and with TRPO with a discount factor of 0.95 (shown to be the best discount factor in the previous experiments). Then for evaluation, we run the trained ATRPO and TRPO policies for a trajectory of 10,000 timesteps (or until the agent falls). We use the same random seeds for the two algorithms. Figure 2 is a plot of the speed of the agent at each time step of the trajectory, using the *seed that*

*gives the best performance for discounted TRPO*. We see in Figure 2 that the discounted algorithm gives a higher initial speed at the beginning of the trajectory. However its overall speed is much more erratic throughout the trajectory, resulting in the agent falling over after approximately 5000 steps. This coincides with the notion of discounting where more emphasis is placed at the beginning of the trajectory and ignores longer-term behavior. On the other hand, the average-reward policy — while having a slightly lower velocity overall throughout its trajectory — is able to sustain the trajectory much longer, thus giving it a higher total return. In fact, we observed that for all 10 random seeds we tested, the average reward agent is able to finish the entire 10,000 time step trajectory without falling. In Table 1 we present the summary statistics of trajectory length for all trajectories using discounted TRPO we note that the median trajectory length for the TRPO discounted agent is 452.5, meaning that on average TRPO performs significantly worse than what is reported in Figure. 2.

Table 1. Summary statistics for all 10 trajectories using a Humanoid-v3 agent trained with TRPO

Min	Max	Average	Median	Std
108	4806	883.1	452.5	1329.902

## 7. Related Work

Dynamic programming algorithms for finding the optimal average reward policies have been well-studied (Howard, 1960; Blackwell, 1962; Veinott, 1966). Several tabular Q-learning-like algorithms for problems with unknown dynamics have been proposed, such as R-Learning (Schwartz, 1993), RVI Q-Learning (Abounadi et al., 2001), CSV-Learning (Yang et al., 2016), and Differential Q-Learning (Wan et al., 2020). Mahadevan (1996) conducted a thorough empirical analysis of the R-Learning algorithm. We note that much of the previous work on average reward RL



focuses on the tabular setting without function approximations, and the theoretical properties of many of these Q-learning-based algorithms are not well understood (in particular R-learning). More recently, POLITEX updates policies using a Boltzmann distribution over the sum of action-value function estimates of the previous policies (Abbasi-Yadkori et al., 2019) and Wei et al. (2020) introduced a model-free algorithm for optimizing the average reward of weakly-communicating MDPs.

For policy gradient methods, Baxter & Bartlett (2001) showed that if  $1/(1 - \gamma)$  is large compared to the mixing time of the Markov chain induced by the MDP, then the gradient of  $\rho_\gamma(\pi)$  can accurately approximate the gradient of  $\rho(\pi)$ . Kakade (2001a) extended upon this result and provided an error bound on using an optimal discounted policy to maximize the average reward. In contrast, our work directly deals with the average reward objective and provides theoretical guidance on the optimal step size for each policy update.

Policy improvement bounds have been extensively explored in the discounted case. The results from Schulman et al. (2015) are extensions of Kakade & Langford (2002). Pirota et al. (2013) also proposed an alternative generalization to Kakade & Langford (2002). Achiam et al. (2017) improved upon Schulman et al. (2015) by replacing the maximum divergence with the average divergence.

## 8. Conclusion

In this paper, we introduce a novel policy improvement bound for the average reward criterion. The bound is based on the average divergence between two policies and Kemeny’s constant or mixing time of the Markov chain. We show that previous existing policy improvement bounds for the discounted case results in a non-meaningful bound for the average reward objective. Our work provides the theoretical justification and the means to generalize the popular trust-region based algorithms to the average reward setting. Based on this theory, we propose ATRPO, a modification of the TRPO algorithm for on-policy DRL. We demonstrate through a series of experiments that ATRPO is highly effective on high-dimensional continuing control tasks.

## Acknowledgements

We would like to extend our gratitude to Quan Vuong and the anonymous reviewers for their constructive comments and suggestions. We also thank Shuyang Ling, Che Wang, Zining (Lily) Wang, and Yanqiu Wu for the insightful discussions on this work.

## References

- Abbasi-Yadkori, Y., Bartlett, P., Bhatia, K., Lazic, N., Szepesvari, C., and Weisz, G. Politex: Regret bounds for policy iteration using expert prediction. In *International Conference on Machine Learning*, pp. 3692–3702, 2019.
- Abounadi, J., Bertsekas, D., and Borkar, V. S. Learning algorithms for markov decision processes with average cost. *SIAM Journal on Control and Optimization*, 40(3): 681–698, 2001.
- Achiam, J., Held, D., Tamar, A., and Abbeel, P. Constrained policy optimization. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 22–31. JMLR. org, 2017.
- Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. Optimality and approximation with policy gradient methods in markov decision processes. In *Conference on Learning Theory*, pp. 64–66. PMLR, 2020.
- Altman, E. *Constrained Markov decision processes*, volume 7. CRC Press, 1999.
- Amit, R., Meir, R., and Ciosek, K. Discount factor as a regularizer in reinforcement learning. In *International conference on machine learning*, 2020.
- Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., and Mané, D. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- Andrychowicz, M., Raichuk, A., Stańczyk, P., Orsini, M., Girgin, S., Marinier, R., Hussenot, L., Geist, M., Pietquin, O., Michalski, M., et al. What matters in on-policy reinforcement learning? a large-scale empirical study. *arXiv preprint arXiv:2006.05990*, 2020.
- Baxter, J. and Bartlett, P. L. Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 15:319–350, 2001.
- Bertsekas, D. P., Bertsekas, D. P., Bertsekas, D. P., and Bertsekas, D. P. *Dynamic programming and optimal control*, volume 1,2. Athena scientific Belmont, MA, 1995.
- Blackwell, D. Discrete dynamic programming. *The Annals of Mathematical Statistics*, pp. 719–726, 1962.
- Boyd, S., Boyd, S. P., and Vandenberghe, L. *Convex optimization*. Cambridge university press, 2004.
- Brémaud, P. *Markov Chains Gibbs Fields, Monte Carlo Simulation and Queues*. Springer, 2 edition, 2020.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym, 2016.

- Cho, G. E. and Meyer, C. D. Comparison of perturbation bounds for the stationary distribution of a markov chain. *Linear Algebra and its Applications*, 335(1-3):137–150, 2001.
- Even-Dar, E., Kakade, S. M., and Mansour, Y. Online markov decision processes. *Mathematics of Operations Research*, 34(3):726–736, 2009.
- Grinstead, C. M. and Snell, J. L. *Introduction to probability*. American Mathematical Soc., 2012.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *International Conference on Machine Learning (ICML)*, 2018.
- Horn, R. A. and Johnson, C. R. *Matrix analysis*. Cambridge university press, 2012.
- Howard, R. A. *Dynamic programming and markov processes*. John Wiley, 1960.
- Hunter, J. J. Stationary distributions and mean first passage times of perturbed markov chains. *Linear Algebra and its Applications*, 410:217–243, 2005.
- Jiang, N., Kulesza, A., Singh, S., and Lewis, R. The dependence of effective planning horizon on model accuracy. In *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*, pp. 1181–1189. Citeseer, 2015.
- Jiang, N., Singh, S. P., and Tewari, A. On structural properties of mdps that bound loss due to shallow planning. In *IJCAI*, pp. 1640–1647, 2016.
- Kakade, S. Optimizing average reward using discounted rewards. In *International Conference on Computational Learning Theory*, pp. 605–615. Springer, 2001a.
- Kakade, S. and Langford, J. Approximately optimal approximate reinforcement learning. In *International Conference on Machine Learning*, volume 2, pp. 267–274, 2002.
- Kakade, S. M. A natural policy gradient. *Advances in neural information processing systems*, 14, 2001b.
- Kallenberg, L. *Linear Programming and Finite Markovian Control Problems*. Centrum Voor Wiskunde en Informatica, 1983.
- Kemeny, J. and Snell, I. *Finite Markov Chains*. Van Nostrand, New Jersey, 1960.
- Lehmann, E. L. and Casella, G. *Theory of point estimation*. Springer Science & Business Media, 2006.
- Lehnert, L., Laroche, R., and van Seijen, H. On value function representation of long horizon problems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Levin, D. A. and Peres, Y. *Markov chains and mixing times*, volume 107. American Mathematical Soc., 2017.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. *International Conference on Learning Representations (ICLR)*, 2016.
- Mahadevan, S. Average reward reinforcement learning: Foundations, algorithms, and empirical results. *Machine learning*, 22(1-3):159–195, 1996.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. Playing atari with deep reinforcement learning. *NIPS Deep Learning Workshop*, 2013.
- Naik, A., Shariff, R., Yasui, N., and Sutton, R. S. Discounted reinforcement learning is not an optimization problem. *NeurIPS Optimization Foundations for Reinforcement Learning Workshop*, 2019.
- Neu, G., Antos, A., György, A., and Szepesvári, C. Online markov decision processes under bandit feedback. In *Advances in Neural Information Processing Systems*, pp. 1804–1812, 2010.
- Peters, J. and Schaal, S. Reinforcement learning of motor skills with policy gradients. *Neural networks*, 21(4):682–697, 2008.
- Petrik, M. and Scherrer, B. Biasing approximate dynamic programming with a lower discount factor. In *Twenty-Second Annual Conference on Neural Information Processing Systems-NIPS 2008*, 2008.
- Pirotta, M., Restelli, M., Pecorino, A., and Calandriello, D. Safe policy iteration. In *International Conference on Machine Learning*, pp. 307–315, 2013.
- Ross, K. W. Constrained markov decision processes with queueing applications. *Dissertation Abstracts International Part B: Science and Engineering*[DISS. ABST. INT. PT. B- SCI. & ENG.], 46(4), 1985.
- 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. *International Conference on Learning Representations (ICLR)*, 2016.

- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Schwartz, A. A reinforcement learning method for maximizing undiscounted rewards. In *Proceedings of the tenth international conference on machine learning*, volume 298, pp. 298–305, 1993.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484, 2016.
- Silver, D., Hubert, T., Schrittwieser, J., Antonoglou, I., Lai, M., Guez, A., Lanctot, M., Sifre, L., Kumaran, D., Graepel, T., et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.
- Song, H. F., Abdolmaleki, A., Springenberg, J. T., Clark, A., Soyer, H., Rae, J. W., Noury, S., Ahuja, A., Liu, S., Tirumala, D., et al. V-mpo: on-policy maximum a posteriori policy optimization for discrete and continuous control. *International Conference on Learning Representations*, 2020.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. 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*, pp. 1057–1063, 2000.
- Tadepalli, P. and Ok, D. H-learning: A reinforcement learning method to optimize undiscounted average reward. Technical Report 94-30-01, Oregon State University, 1994.
- Tessler, C., Mankowitz, D. J., and Mannor, S. Reward constrained policy optimization. *International Conference on Learning Representation (ICLR)*, 2019.
- Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033. IEEE, 2012.
- Tsybakov, A. B. *Introduction to nonparametric estimation*. Springer Science & Business Media, 2008.
- Veinott, A. F. On finding optimal policies in discrete dynamic programming with no discounting. *The Annals of Mathematical Statistics*, 37(5):1284–1294, 1966.
- Vuong, Q., Zhang, Y., and Ross, K. W. Supervised policy update for deep reinforcement learning. In *International Conference on Learning Representation (ICLR)*, 2019.
- Wan, Y., Naik, A., and Sutton, R. S. Learning and planning in average-reward markov decision processes. *arXiv preprint arXiv:2006.16318*, 2020.
- Wei, C.-Y., Jafarnia-Jahromi, M., Luo, H., Sharma, H., and Jain, R. Model-free reinforcement learning in infinite-horizon average-reward markov decision processes. In *International conference on machine learning*, 2020.
- Wu, Y., Mansimov, E., Grosse, R. B., Liao, S., and Ba, J. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. In *Advances in neural information processing systems (NIPS)*, pp. 5285–5294, 2017.
- Yang, S., Gao, Y., An, B., Wang, H., and Chen, X. Efficient average reward reinforcement learning using constant shifting values. In *AAAI*, pp. 2258–2264, 2016.
- Yang, T.-Y., Rosca, J., Narasimhan, K., and Ramadge, P. J. Projection-based constrained policy optimization. In *International Conference on Learning Representation (ICLR)*, 2020.
- Zhang, Y., Vuong, Q., and Ross, K. First order constrained optimization in policy space. *Advances in Neural Information Processing Systems*, 33, 2020.
- Zhao, T., Hachiya, H., Niu, G., and Sugiyama, M. Analysis and improvement of policy gradient estimation. In *Advances in Neural Information Processing Systems*, pp. 262–270, 2011.