# On the Convergence of Reinforcement Learning with Monte Carlo Exploring Starts $^\star$

## Jun Liu

*Department of Applied Mathematics*
*University of Waterloo*
*Waterloo, Ontario N2L 3G1, Canada*

**Abstract**

A basic simulation-based reinforcement learning algorithm is the Monte Carlo Exploring States (MCES) method, also known as optimistic policy iteration, in which the value function is approximated by simulated returns and a greedy policy is selected at each iteration. The convergence of this algorithm in the general setting has been an open question. In this paper, we investigate the convergence of this algorithm for the case with undiscounted costs, also known as the stochastic shortest path problem. The results complement existing partial results on this topic and thereby helps further settle the open problem. As a side result, we also provide a proof of a version of the supermartingale convergence theorem commonly used in stochastic approximation.

*Key words:* Reinforcement Learning; Markov Decision Processes; Stochastic Control; Monte Carlo Exploring States; Optimistic Policy Iteration; Convergence; Stochastic Shortest Path Problem.

## 1 Introduction

Reinforcement learning has gained tremendous popularity in recent years [11]. Simulation-based methods for reinforcement learning or stochastic control have achieved notable success [10]. One particularly simple simulation-based method, called Monte Carlo Exploring Starts (MCES), was introduced in detail in the classic book by Sutton and Barto [11, Chapter 5]. In this method, the value function is estimated by the average simulated returns and the policy is updated using a greedy policy based on the current estimate of the value function. Because of its fundamental simplicity and importance, Sutton and Barto stated that the convergence of the MCES algorithm to the actual optimal value is "one of the most fundamental open theoretical questions in reinforcement learning" [11, p. 99].

Partial results on convergence analysis of MCES exist in the literature. Most notably, Tsitsiklis [12] proved that MCES, which he termed optimistic policy iteration, converges under two assumptions. First, each state is selected for updating with the same frequency. Second, the problem is strictly discounted with a discount factor less

than one. This result was extended to the undiscounted case by Chen [5] under the assumption that all policies are proper (i.e., reaching a terminal state is inevitable under all policies). A more recent result by Wang and Ross [13] proved convergence of MCES under the assumption of optimal policy feed-forward environments, where states cannot be revisited under an optimal policy. We note that the approach taken in [13] mostly uses finite graph and probabilistic argument, whereas the approach in [12] (and also [5]) is along the lines of stochastic approximation [4,9].

In this paper, we investigate the convergence of MCES/optimistic policy iteration in the undiscounted case without the assumption of optimal policy feed-forward environments and without the assumption that all polices are proper. Compared with the results in [12,5,13], we consider both uniform and nonuniform exploring starts. In the uniform case, we extend the results of Tsitsiklis [12] to the undiscounted, i.e., stochastic shortest path problem. For the case that all policies are proper, our proof differs from that in [5] and is closer in spirit to that of [12] (see how Lemma 7 generalizes Lemma 2 in [12]). We also discuss how to work around the proper policy assumption. In the nonuniform case, we argue that the choice of stepsize should be component-dependent to agree with the clas-

---

$^\star$ This paper was not presented at any IFAC meeting.
*Email address:* `j.liu@uwaterloo.ca` (Jun Liu).

sical version of MCES discussed in [11]. We believe the convergence results established here could help further settle the long-standing open problem. As a side result, we also prove a version of the supermartingale convergence theorem that is commonly used in the literature of stochastic approximation, whose proof, however, is not available in classic books such as [4]. Furthermore, we provide an alternative and hopefully more direct treatment of stochastic approximations directly based on the supermartingale martingale convergence theorem (cf. Chapter 4 of the classic book [4]), which may be of independent interest.

The paper is organized as follows. In Section 2, we present the problem formulation and the preliminaries for proving the convergence of MCES/optimistic policy iteration. In Section 3, we present the convergence proof for the case that all policies are proper. We discuss the case without the proper policy assumption in Section 4 and the case with nonuniform initial exploring in Section 5. A simple illustrative example is presented in Section 6. Some concluding remarks are presented in Section 7. The Appendix includes a self-contained treatment of supermartingale convergence and stochastic approximation results.

## 2 Problem formulation and preliminaries

### 2.1 Markov decision problem

Let $M = (S, A, P)$ be a *Markov decision process*, where $S = \{1, \cdots, n\}$ is a finite set of states, $A$ is a finite set of actions, and $P : S \times A \times S \to [0, 1]$ is a transition probability function. For each action $a \in A$, we can represent $P(\cdot, a, \cdot)$ as a matrix $P(a)$ whose entries $P_{ij}(a)$ satisfy

$$P_{ij}(a) = P(i, a, j) = \mathbb{P}(s_{t+1} = j | s_t = i, a_t = a),$$

where $\{(s_t, a_t)_{t=0}^{t=\infty}\} \subseteq S \times A$ is an evolution of the MDP $M$. In words, $P_{ij}(a)$ denotes the probability of having a transition from the state $i$ to the state $j$ under the action $a$.

A *policy* is a function $\mu : S \to A$. Clearly, the set of all policies is finite. We denote this set by $\Pi$. Given a policy $\mu \in \Pi$, we define the cost-to-go value [1] of the policy starting from a state $i$ as

$$J^{\mu}(i) = \mathbb{E}\left[\sum_{t=0}^{\infty} \alpha^t g(s_t, \mu(s_t)) | s_0 = i\right],$$

---

[1] We use a cost function formulation as commonly seen in stochastic control, which is equivalent to a reward function formulation in reinforcement learning, albeit the difference of using minimization in place of maximization for values.

where $\{s_t\}_{t=0}^{\infty}$ is a state evolution under the policy $\mu$, $g : S \times A \to \mathbb{R}$ is the stage cost, and $\alpha \in [0, 1]$ is a discount factor. The optimal cost-to-go value $J^*$ is defined as

$$J^*(i) = \min_{\mu \in \Pi} J^{\mu}(i).$$

Since the set of policies is finite, the optimal value is always attainable by an optimal policy. That is, there exists $\mu^* \in \Pi$ such that $J^{\mu^*} = J^*$. A Makov decision problem often is concerned with finding the optimal value $J^*$ and an optimal policy $\mu^*$ (which may not be unique).

We will primarily be focusing on the so-called *stochastic shortest path problem* in this paper, i.e. the Markov decision problem above with $\alpha = 1$. To make the cost-to-go value well-defined, we assume that there exists a terminal state, denoted by 0, and modify the transition probability function to satisfy $\sum_{j=1}^{n} P_{ij}(a) \leq 1$ and $P_{i0}(a) = 1 - \sum_{j=1}^{n} P_{ij}(a)$ for all $i \in S$ and $a \in A$. In addition, the terminal state is assumed to be a trap state in the sense that $P_{00}(a) = 1$ and $P_{0j}(a) = 0$ for all $a \in A$ and $j \in S$. We also assume $g(0, a) = 0$ for all $a \in A$ such that $J^{\mu}(0) = 0$ for all $\mu \in \Pi$. Hence we do not need to discuss the value at state 0.

### 2.2 Dynamic programming operators

We define two dynamic programming operators $T_{\mu} : \mathbb{R}^n \to \mathbb{R}^n$ and $T : \mathbb{R}^n \to \mathbb{R}^n$ as follows. Given $J \in \mathbb{R}^n$ and $\mu \in \Pi$, let

$$T^{\mu} J(i) = g(i, \mu(i)) + \alpha \sum_{j=1}^{n} P_{ij}(\mu(i)) J(j), \qquad (1)$$

and

$$T J(i) = \min_{a \in A}\left[g(i, a) + \alpha \sum_{j=1}^{n} P_{ij}(a) J(j)\right]. \qquad (2)$$

For convenience, we can write (1) in a vector format as

$$T_{\mu} J = g_{\mu} + \alpha P_{\mu} J,$$

where $g_{\mu} = [g(1, \mu(1))\ g(2, \mu(2))\ \cdots\ g(n, \mu(n))]^T \in \mathbb{R}^n$ and $P_{\mu} = (P_{ij}(\mu(i))) \in \mathbb{R}^{n \times n}$. It follows that, for each $J \in \mathbb{R}^n$, there exists $\mu \in \Pi$ such that

$$T J = T_{\mu} J.$$

Such a policy is called a *greedy policy* corresponding to $J$.

### 2.3 Optimistic policy iteration with Monte Carlo policy evaluation

Following [12], we can write the main procedure of optimistic policy iteration using Monte Carlo simulations

for policy evaluation as

$$J_{t+1} = (1 - \gamma_t)J_t + \gamma_t(J^{\mu_t} + w_t), \qquad (3)$$

where $J_t$ is the current value vector, $\gamma_t$ is a scalar stepsize parameter (time-varying but deterministic), and $J^{\mu_t}$ is the expected cost value of the current policy $\mu_t$. Given the current value $J_t$, a greedy policy $\mu_t$ is chosen according to

$$T_{\mu_t}J_t = TJ_t. \qquad (4)$$

The noise $w_t$ captures the discrepancy between the expected cost $J^{\mu_t}$ and observed cumulative cost $J^{\mu_t} + w_t$. Let $\mathcal{F}_t$ be the natural filtration generated by the process (3). Since the observed cumulative cost gives an unbiased estimate $J^{\mu_t}$, we have $\mathbb{E}[w_t|\mathcal{F}_t] = 0$. Furthermore, the variance of $w_t$ (conditioned on $\mathcal{F}_t$) is only a function of the initial state and the current policy $\mu_t$. Because the numbers of states and polices are finite, we also have $\mathbb{E}[\|w_t\|^2 | \mathcal{F}_t] \leq C$, for some constant $C$.

## 2.4 Preliminaries

We present some technical preliminaries for convergence analysis. We focus on the shortest path problem (i.e. $\alpha = 1$). A policy $\mu \in \Pi$ is said to be proper if the terminal state 0 is reached with probability 1 from any initial state.

**Assumption 1** *All policies in $\Pi$ are proper.*

**Assumption 2** *The stepsize parameter satisfies $\sum_{t=0}^{\infty} \gamma_t = \infty$ and $\sum_{t=0}^{\infty} \gamma_t^2 < \infty$.*

Based on Assumption 1, a well-known result is that the dynamic programming operators $T$ and $T_\mu$ are contractive with respect to a weighted maximum norm.

**Lemma 3** *[4, Proposition 2.2, p. 23] If Assumption 1 holds, then there exists some $\beta \in [0,1)$ and a vector $\theta \in \mathbb{R}^n$ of positive components such that*

$$\sum_{j=1}^{n} P_{ij}(a)\theta(j) \leq \beta\theta(i), \quad \forall i \in S, \quad \forall a \in A.$$

*In particular, this statement implies that*

$$\|T_\mu J_1 - T_\mu J_2\|_\theta \leq \beta \|J_1 - J_2\|_\theta, \quad \forall \mu \in \Pi, \forall J_1, J_2 \in \mathbb{R}^n,$$

*and*

$$\|TJ_1 - TJ_2\|_\theta \leq \beta \|J_1 - J_2\|_\theta, \quad \forall J_1, J_2 \in \mathbb{R}^n,$$

*where the weighted maximum norm $\|\cdot\|_\theta$ is defined by $\|J\|_\theta = \max_{1 \leq i \leq n} \frac{|J(i)|}{\theta(i)}$.*

Let $\theta \in \mathbb{R}^n$ be a vector of positive components. Define $\Theta = \text{diag}\{\theta(1), \theta(2), \cdots, \theta(n)\}$. Let $\mathbf{1} \in \mathbb{R}^n$ be the column vector with all components equal to 1. The above lemma shows that, in matrix form,

$$P_\mu \Theta \mathbf{1} \leq \beta \Theta \mathbf{1}, \quad \forall \mu \in \Pi, \qquad (5)$$

where the inequality is interpreted component-wise[2]. We refer to this as a weighted contractive property for $P_\mu$.

We also recall the following property on the dynamical programming parameters $T$ and $T_\mu$ for a stochastic shortest path problem.

**Lemma 4** *[4, Lemma 2.2, p. 21] For every scalar $c \geq 0$, $J \in \mathbb{R}^n$, and $\mu \in \Pi$, we have*

$$T(J + c\mathbf{1}) \leq TJ + c\mathbf{1}, \quad T_\mu(J + c\mathbf{1}) \leq T_\mu J + c\mathbf{1}, \quad (6)$$

*where $c$ is any nonnegative scalar. If $c$ is negative, then the inequalities are reversed.*

We can also prove a slight modification of the above lemma using (5).

**Lemma 5** *Suppose that Assumption 1 holds. For every scalar $c \geq 0$, $J \in \mathbb{R}^n$, and $\mu \in \Pi$, we have*

$$T(J + c\Theta\mathbf{1}) \leq TJ + \beta c\Theta\mathbf{1}, \quad T_\mu(J + c\Theta\mathbf{1}) \leq T_\mu J + \beta c\Theta\mathbf{1},$$

*where $c$ is any nonnegative scalar. If $c$ is negative, then the inequalities are reversed.*

**Proof** Let $\mu$ be a greedy policy corresponding to $J$, i.e., $T_\mu J = TJ$. By (5), we have

$$T(J + c\Theta\mathbf{1}) \leq T_\mu(J + c\Theta\mathbf{1}) = g_\mu + P_\mu(J + c\Theta\mathbf{1})$$
$$= g_\mu + P_\mu J + P_\mu c\Theta\mathbf{1} = T_\mu J + P_\mu c\Theta\mathbf{1}$$
$$\leq T_\mu J + \beta c\Theta\mathbf{1} = TJ + \beta c\Theta\mathbf{1}.$$

The above also shows $T_\mu(J + c\Theta\mathbf{1}) \leq T_\mu J + \beta c\Theta\mathbf{1}$ for any $\mu \in \Pi$. ∎

By the contraction mapping theorem, Lemma 3 implies the following convergence result.

**Lemma 6** *[6,4] If Assumption 1 holds, we have, for every $J \in \mathbb{R}^n$ and $\mu \in \Pi$,*

$$\lim_{t \to \infty} T^t J = J^*, \quad \lim_{t \to \infty} T_\mu^t J = J^\mu,$$

*where $J^*$ and $J^\mu$ are the unique fixed points of $T$ and $T_\mu$, respectively.*

---

[2] In the sequel, all vector inequalities are interpreted component-wise.

The next lemma is a modified version of Lemma 2 in [12]. Let $\theta \in \mathbb{R}^n$ and $\Theta = \mathrm{diag}\{\theta(1), \theta(2), \cdots, \theta(n)\}$ be defined above. For the sequence $\{J_t\}_{t=0}^{\infty} \subset \mathbb{R}^n$, define

$$c_t = TJ_t - J_t, \quad \lambda_t = \max(c_t, 0), \quad t \geq 0,$$

where max is taken component-wise. Then clearly $\lambda_t$ is a nonnegative vector and $c_t \leq \lambda_t$.

**Lemma 7** *Suppose that Assumption 1 holds. For every $t \geq 0$, we have*

*(1)* $T_{\mu_t}^k J_t \leq J_t + \frac{\left\| \Theta^{-1} \lambda_t \right\|_{\infty} \Theta \mathbf{1}}{1-\beta}$, *for all $k \geq 1$,*

*(2)* $J^{\mu_t} \leq J_t + \frac{\left\| \Theta^{-1} \lambda_t \right\|_{\infty} \Theta \mathbf{1}}{1-\beta}$,

*(3)* $J^{\mu_t} \leq TJ_t + \frac{\beta \left\| \Theta^{-1} \lambda_t \right\|_{\infty} \Theta \mathbf{1}}{1-\beta}$.

**Proof** Note that $\left\| \Theta^{-1} \lambda_t \right\|_{\infty}$ is the weighted maximum norm of $\lambda_t$ with respect to the vector $\theta$. From (4), we have $T_{\mu_t} J_t = TJ_t$. It follows that $T_{\mu_t} J_t = J_t + c_t$. Applying $T_{\mu_t}$ to both sides of this equation gives

$$T_{\mu_t}^2 J_t = T_{\mu_t}(J_t + c_t) = g_{\mu_t} + P_{\mu_t}(J_t + c_t)$$
$$= T_{\mu_t} J_t + P_{\mu_t} c_t = J_t + c_t + P_{\mu_t} c_t.$$

By induction, we obtain

$$T_{\mu_t}^k J_t = J_t + (I + P_{\mu_t} + P_{\mu_t}^2 + \cdots + P_{\mu_t}^{k-1}) c_t. \quad (7)$$

We have, for $m \geq 1$,

$$P_{\mu_t}^m c_t = P_{\mu_t}^m \Theta \Theta^{-1} c_t \leq P_{\mu_t}^m \Theta \Theta^{-1} \lambda_t \leq P_{\mu_t}^m \Theta \left\| \Theta^{-1} \lambda_t \right\|_{\infty} \mathbf{1}$$
$$= \left\| \Theta^{-1} \lambda_t \right\|_{\infty} P_{\mu_t}^{m-1} P_{\mu_t} \Theta \mathbf{1} \leq \left\| \Theta^{-1} \lambda_t \right\|_{\infty} P_{\mu_t}^{m-1} \beta \Theta \mathbf{1}$$
$$= \beta \left\| \Theta^{-1} \lambda_t \right\|_{\infty} P_{\mu_t}^{m-1} \Theta \mathbf{1}$$
$$\leq \beta^m \left\| \Theta^{-1} \lambda_t \right\|_{\infty} \Theta \mathbf{1}, \quad (8)$$

where the first two inequalities follow from the fact that elements of $P_{\mu_t}^{m-1} P_{\mu_t} \Theta$ and $P_{\mu_t}^{m-1} P_{\mu_t} \Theta \Theta^{-1}$ are nonnegative and we can bound components of $c_t$ with $\lambda_t$ and $\Theta^{-1} \lambda_t$ with $\left\| \Theta^{-1} \lambda_t \right\|_{\infty} \mathbf{1}$, the third inequality follows from (5), the last inequality follows from an inductive argument, and the equations follow from straightforward rearrangements. Part of the above inequality also shows that, for $m = 0$, $c_t \leq \left\| \Theta^{-1} \lambda_t \right\|_{\infty} \Theta \mathbf{1}$. Hence by (7) we obtain

$$T_{\mu_t}^k J_t \leq J_t + (1 + \beta + \beta^2 + \cdots + \beta^{k-1}) \left\| \Theta^{-1} \lambda_t \right\|_{\infty} \Theta \mathbf{1}$$
$$\leq J_t + \frac{\left\| \Theta^{-1} \lambda_t \right\|_{\infty} \Theta \mathbf{1}}{1 - \beta}, \quad (9)$$

where in the first inequality we used (8) and the fact that $c_t \leq \lambda_t \left\| \Theta^{-1} \lambda_t \right\|_{\infty} \Theta \mathbf{1}$. We proved item (1). Since $\lim_{k \to \infty} T_{\mu_t}^k J_t = J^{\mu_t}$, we proved item (2) by letting $k \to$

$\infty$. Finally, applying $T_{\mu_t}$ to both sides of the inequality in item (1) and using the fact $J^{\mu_t} = T_{\mu_t} J^{\mu_t}$, we obtain

$$J^{\mu_t} = T_{\mu_t} J^{\mu_t} \leq T_{\mu_t}(J_t + \frac{\left\| \Theta^{-1} \lambda_t \right\|_{\infty} \Theta \mathbf{1}}{1 - \beta}).$$

By Lemma 5 and the fact that $T_{\mu_t} J_t = TJ_t$, we obtained item (3). ∎

## 3 Convergence analysis for the stochastic shortest path problem with proper policies

The convergence analysis starts with an asymptotic estimate for $c_t = TJ_t - J_t$ and $\lambda_t = \max(c_t, 0)$. All convergence and asymptotic estimates for random variables in this section are understood in the sense of probability 1.

**Lemma 8** *[12] Under Assumption 2, we have*

$$\limsup_{t \to \infty} c_t \leq 0 \ and \ \lim_{t \to \infty} \lambda_t = 0.$$

**Proof** The proof for $\limsup_{t \to \infty} c_t \leq 0$ was established in [12] for the case $\alpha < 1$. The same argument holds for $\alpha = 1$. Here is an outline of the proof. Since $T_{\mu_t} J = g_{\mu_t} + P_{\mu_t} J$ for any $J \in \mathbb{R}^n$, we can verify that

$$TJ_{t+1} \leq T_{\mu_t} J_{t+1} = T_{\mu_t}((1 - \gamma_t)J_t + \gamma_t J^{\mu_t} + \gamma_t w_t)$$
$$= J_{t+1} + (1 - \gamma_t)(TJ_t - J_t) + \gamma_t v_t,$$

where we need to use the fact that $TJ_t = T_{\mu_t} J_t$ and $v_t = P_{\mu_t} w_t - w_t$. By the property on $w_t$, we have $\mathbb{E}[v_t \mid \mathcal{F}_t] = 0$ and $\mathbb{E}[\|v_t\|^2 \mid \mathcal{F}_t] \leq C'$ for some constant $C'$. Hence, $c_t$ satisfies

$$c_{t+1} \leq (1 - \gamma_t)c_t + \gamma_t v_t.$$

Consider another iteration

$$V_{t+1} = (1 - \gamma_t)V_t + \gamma_t v_t.$$

If $V_0 = c_0$, then a comparison argument shows that $c_t \leq V_t$ for all $t \geq 0$. By a standard supermartingale convergence argument on stochastic iterations [4, Chapter 4, p. 143] (see also Proposition 23 and Lemma 24 in the Appendix), one can show that $V_t$ converges to 0 in probability 1. Hence, $\limsup_{t \to \infty} c_t \leq 0$. Since $\lambda_t = \max(c_t, 0)$, it follows that $\lim_{t \to \infty} \lambda_t = 0$. ∎

In particular, the above lemma shows that, for any $\varepsilon > 0$, there exists $t(\varepsilon)$ such that

$$\frac{\beta \left\| \Theta^{-1} \lambda_t \right\|_{\infty} \Theta \mathbf{1}}{1 - \beta} \leq \varepsilon \Theta \mathbf{1}, \quad \forall t \geq t(\varepsilon).$$

Putting this into Lemma 7(3) shows that

$$J_{\mu_t} \le TJ_t + \varepsilon\Theta\mathbf{1}, \quad \forall t \ge t(\varepsilon).$$

By (3), we obtain

$$
\begin{aligned}
J_{t+1} &= (1 - \gamma_t)J_t + \gamma_t(J^{\mu_t} + w_t) \\
&\le (1 - \gamma_t)J_t + \gamma_t TJ_t + \gamma_t\varepsilon\Theta\mathbf{1} + \gamma_t w_t, \quad \forall t \ge t(\varepsilon).
\end{aligned}
$$

Define a mapping $H_\varepsilon : \mathbb{R}^n \to \mathbb{R}^n$ as $H_\varepsilon J = TJ + \varepsilon\Theta\mathbf{1}$ and consider the sequence $\{Z_t\}$ generated by

$$
\begin{aligned}
Z_{t+1} &= (1 - \gamma_t)Z_t + \gamma_t TZ_t + \gamma_t\varepsilon\Theta\mathbf{1} + \gamma_t w_t \\
&= (1 - \gamma_t)Z_t + \gamma_t(H_\varepsilon Z_t + w_t), \quad t \ge t(\varepsilon),
\end{aligned}
$$

and $Z_{t(\varepsilon)} = J_{t(\varepsilon)}$. Then, by comparison,

$$J_t \le Z_t, \quad \forall t \ge t(\varepsilon). \tag{10}$$

Since $T$ is a contraction under the weighted maximum norm $\|\cdot\|_\theta$, so is $H_\varepsilon$. By Proposition 4.4 in [4] (see also Proposition 23 in the Appendix), we know that $Z_t$ converges to the unique fixed point of $H_\varepsilon$, denoted by $Z_\varepsilon^*$.

The following lemma estimates the fixed point of $H_\varepsilon$ relative to $J^*$.

**Lemma 9** *Under Assumption 1, we have*

$$J^* - \frac{\varepsilon}{1 - \beta}\Theta\mathbf{1} \le Z_\varepsilon^* \le J^* + \frac{\varepsilon}{1 - \beta}\Theta\mathbf{1}.$$

**Proof** By Lemma 5 and $TJ^* = J^*$, we have

$$
\begin{aligned}
H_\varepsilon(J^* + \frac{\varepsilon}{1 - \beta}\Theta\mathbf{1}) &= T(J^* + \frac{\varepsilon}{1 - \beta}\Theta\mathbf{1}) + \varepsilon\Theta\mathbf{1} \\
&\le TJ^* + \frac{\varepsilon\beta}{1 - \beta}\Theta\mathbf{1} + \varepsilon\Theta\mathbf{1} \\
&= J^* + \frac{\varepsilon}{1 - \beta}\Theta\mathbf{1}.
\end{aligned}
$$

It follows that

$$Z_\varepsilon^* = \lim_{k \to \infty} H_\varepsilon^k(J^* + \frac{\varepsilon}{1 - \beta}\Theta\mathbf{1}) \le J^* + \frac{\varepsilon}{1 - \beta}\Theta\mathbf{1},$$

where we used monotonicity of $H_\varepsilon$ (implied by that of $T$). Similary, by Lemma 5, we can show that

$$H_\varepsilon(J^* - \frac{\varepsilon}{1 - \beta}\Theta\mathbf{1}) \ge J^* - \frac{\varepsilon}{1 - \beta}\Theta\mathbf{1}$$

and $Z_\varepsilon^* \ge J^* - \frac{\varepsilon}{1-\beta}\Theta\mathbf{1}$. ∎

We now state and prove the main result of the paper.

**Theorem 10** *Under Assumptions 1 an 2, the sequence $J_t$ generated by the optimistic policy iteration (3) and (4), applied to a stochastic shortest path problem, converges to $J^*$, with probability 1.*

**Proof** Given any $\varepsilon > 0$, by the argument preceding (10), there exists $t(\varepsilon)$ such that $J_t \le Z_t$ for all $t \ge t(\varepsilon)$. Since $\lim_{t \to \infty} Z_t = Z_\varepsilon^*$, it follows that $\limsup_{t \to \infty} J_t \le Z_\varepsilon^*$. By Lemma 9, we have $\limsup_{t \to \infty} J_t \le J^* + \frac{\varepsilon}{1 - \beta}\Theta\mathbf{1}$. Since the choice of $\varepsilon > 0$ is arbitrary, we obtain $\limsup_{t \to \infty} J_t \le J^*$. By the definition of $J^{\mu_t}$ and $J^*$, we have $J^{\mu_t} \ge J^*$. Hence, (3) implies

$$J_{t+1} \ge (1 - \gamma_t)J_t + \gamma_t J^* + \gamma_t w_t.$$

Consider the iteration

$$Y_{t+1} = (1 - \gamma_t)Y_t + \gamma_t J^* + \gamma_t w_t$$

with $Y_0 = J_0$. Then the sequence $\{Y_t\}$ converges to $J^*$ (see Proposition 4.4 in [4] or Proposition 23 in the Appendix). By comparison, $\liminf_{t \to \infty} J_t \ge J^*$. Hence, $\lim_{t \to \infty} J_t = J^*$. ∎

## 4 Relaxing the proper policy assumption

Assumption 1 requires that all policies are proper. In this section, we discuss how to relax this assumption. For the stochastic shortest path problem, the following relaxed assumption was proposed in [3] (see also [4, Chapter 2]).

**Assumption 11** *There exists at least one proper policy, and every improper policy yields an infinite cost for at least one initial state, i.e., for every improper $\mu \in \Pi$,*

$$J^\mu(i) = \lim_{k \to \infty}[\sum_{t=0}^{k-1} P_\mu g_\mu]_i = \infty \text{ for some } i \in S.$$

**Lemma 12** *[3] If Assumption 11 holds, we have, for every $J \in \mathbb{R}^n$ and proper $\mu \in \Pi$,*

$$\lim_{t \to \infty} T^t J = J^*, \quad \lim_{t \to \infty} T_\mu^t J = J^\mu,$$

*where $J^*$ and $J^\mu$ are the unique fixed points of $T$ and $T_\mu$, respectively.*

There is a problem, however, to analyze the convergence of the optimistic policy iteration (3) and (4) under Assumption 11. Unlike in the standard policy iteration, we cannot guarantee the greedy policy generated by the optimistic iteration is always proper. Hence the value iteration (3) will possibly attain infinity and become invalid. To overcome this issue, a natural way would be

to let the process terminates with a small probability at each stage. This is equivalent to modifying the Markov decision process $M$, by adding a transition with a small probability to the terminal state under each action, so that it satisfies Assumption 1.

A natural question to ask is whether the optimal value of the modified problem stays close to that of the original problem and whether an optimal policy obtained for the modified problem remains an optimal policy for the original problem.

Formally, we define a modified MDP $\hat{M} = (S, A, \hat{P})$ from the original MDP $M = (S, A, P)$ as follows. Let $\hat{P}_{i0}(a) = P_{i0}(a) + p_\varepsilon \sum_{j=1}^n P_{ij}(a)$ and $\hat{P}_{ij}(a) = (1 - p_\varepsilon)P_{ij}(a)$ for all $i, j \in S$ and $a \in A$, where $p_\varepsilon \in (0, 1)$ is a small probability to be chosen. Then $\hat{P}_\mu = (1 - p_\varepsilon)P_\mu$ for all $\mu \in \Pi$ and $\hat{M}$ satisfies Assumption 1. Let $\hat{J}^*$ denote the optimal value for $\hat{M}$ and $\hat{\mu}^*$ a corresponding optimal policy.

**Proposition 13** *Suppose that $M$ satisfies Assumption 11. For every $\varepsilon > 0$, there exists some $\delta > 0$ such that, if $p_\varepsilon \in (0, \delta)$, then $\left\| \hat{J}^* - J^* \right\| \le \varepsilon$. Furthermore, if $\delta > 0$ is sufficiently small, then $p_\varepsilon \in (0, \delta)$ implies that $\hat{\mu}^*$ is also an optimal policy for $M$.*

**Proof** The proof consists of two main parts. First we show that, by Assumption 11, any improper policy for $M$ necessarily has large cost-to-go value for at least one component and hence cannot be optimal for $\hat{M}$ (even if it becomes proper with the modification), provided that $p_\varepsilon$ is chosen sufficiently small. We then show that the value of a proper policy in $\hat{M}$ remains close to its value in $M$, provided that $p_\varepsilon$ is sufficiently small. As a result, the optimal value remains close and optimal policy remains the same for $p_\varepsilon$ chosen sufficiently small.

Let $\mu$ be an improper policy for $M$. Consider the Jordan normal form of $P_\mu$:

$$Q^{-1}P_\mu Q = \begin{bmatrix} I_p & 0 \\ 0 & C \end{bmatrix},$$

where $Q$ is a nonsingular matrix and $C$ has spectral radius $\rho(C) < 1$. We obtain $I_p$ because the eigenvalue 1 of $P_\mu$ is semisimple [7, p. 696]. The dimension of $I_p$ cannot be zero because otherwise $\mu$ would be a proper policy. It follows that

$$Q^{-1} \sum_{t=0}^{k-1} P_\mu^t Q = \begin{bmatrix} kI_p & 0 \\ 0 & \sum_{t=0}^{k-1} C^t \end{bmatrix}, \qquad (11)$$

where $\sum_{t=0}^{\infty} C^t = (I - C)^{-1}$. Since $\hat{P}_\mu = (1 - p_\varepsilon)P_\mu$, we have

$$Q^{-1} \sum_{t=0}^{k-1} \hat{P}_\mu^t Q = \begin{bmatrix} \frac{1-(1-p_\varepsilon)^k}{p_\varepsilon} I_p & 0 \\ 0 & \sum_{t=0}^{k-1}(1 - p_\varepsilon)^t C^t \end{bmatrix},$$

where $\sum_{t=0}^{\infty}(1 - p_\varepsilon)^t C^t = (I - (1 - p_\varepsilon)C)^{-1}$.

By Assumption 11, $\lim_{k\to\infty} \left[ \sum_{t=0}^{k-1} P_\mu^t g_\mu \right]_i = \infty$ for some $i \in S$. This is equivalent to

$$\lim_{k\to\infty} \left[ Q \begin{bmatrix} kI_p & 0 \\ 0 & \sum_{t=0}^{k-1} C^t \end{bmatrix} Q^{-1} g_\mu \right]_i = \infty,$$

which is again equivalent to

$$\lim_{k\to\infty} \left[ Q \begin{bmatrix} kI_p & 0 \\ 0 & 0 \end{bmatrix} Q^{-1} g_\mu \right]_i = \infty$$

in view of $\sum_{t=0}^{\infty} C^t < \infty$. Since $\frac{1-(1-p_\varepsilon)^k}{p_\varepsilon} \to k$, as $p_\varepsilon \to 0$, and $\sum_{t=0}^{\infty}(1 - p_\varepsilon)^t C^t = (I - (1 - p_\varepsilon)C)^{-1}$ is continuous w.r.t. $p_\varepsilon$ and hence bounded for $p_\varepsilon \in [0, 1]$, it is straightforward to verify that, for any $c > 0$, there exists $\delta > 0$, such that

$$\lim_{k\to\infty} \left[ \sum_{t=0}^{k-1} \hat{P}_\mu^t g_\mu \right]_i$$
$$= \lim_{k\to\infty} \left[ Q \begin{bmatrix} \frac{1-(1-p_\varepsilon)^k}{p_\varepsilon} I_p & 0 \\ 0 & \sum_{t=0}^{k-1}(1 - p_\varepsilon)^t C^t \end{bmatrix} Q^{-1} g_\mu \right]_i$$
$$> c, \quad \forall p_\varepsilon \in (0, \delta]. \qquad (12)$$

Now consider a proper policy $\mu$ for $M$ and let $\Pi_0$ denote the set of all proper policies for $M$. Clearly $\mu$ remains a proper policy for $\hat{M}$. The cost vector for $\mu$ in $\hat{M}$ is the unique solution to

$$\hat{J}^\mu = \hat{T}_\mu \hat{J}^\mu.$$

Note that $\hat{T}_\mu$ changes continuous with respect to $p_\varepsilon$ and $\hat{T}_\mu = T_\mu$ when $p_\varepsilon = 0$. Since $T_\mu$ has a unique solution $J^\mu$, it follows that $\hat{J}^\mu$ also changes continuously with respect to $p_\varepsilon$. Hence for any $\rho > 0$, there exists $\delta > 0$ such that

$$\left\| \hat{J}^\mu - J^\mu \right\| < \rho, \quad \forall p_\varepsilon \in (0, \delta], \forall \mu \in \Pi_0, \qquad (13)$$

because the number of policies is finite. Let $J^*$ be the optimal value of $M$ and define

$$d = \min_{\substack{\mu \in \Pi_0 \\ J^\mu \ne J^*}} \left\| J^\mu - J^* \right\|.$$

6

Choose any $\rho < \frac{d}{2}$ and $\delta$ accordingly such that (13) holds. Choose $c = \max_{\mu \in \Pi_0} J^\mu + \rho$ and reduce $\delta$ accordingly such that (12) holds. In view of (13) and the definition of $c$, any improper policy (w.r.t. $M$) cannot be optimal for $\hat{M}$ for all $p_\varepsilon \in (0, \delta]$. Furthermore, suppose that $\hat{J}^*$ is the optimal value and $\hat{\mu}^*$ is an optimal policy for $\hat{M}$. Then $\hat{\mu}^*$ is proper w.r.t. $M$. We claim that $J^{\hat{\mu}^*} = J^*$ and hence $\hat{\mu}^*$ is an optimal policy for $M$. Suppose this is not the case. Then $J^{\hat{\mu}^*} - J^* \geq d$. Since $\left\| \hat{J}^{\hat{\mu}^*} - J^{\hat{\mu}^*} \right\| < \frac{d}{2}$, it follows that $\hat{J}^{\hat{\mu}^*} > J^* + \frac{d}{2}$. Let $\mu^*$ be a proper optimal policy for $M$. Then (13) implies that $\hat{J}^{\mu^*} < J^{\mu^*} + \frac{d}{2} = J^* + \frac{d}{2}$. Hence $\hat{J}^{\mu^*} < \hat{J}^{\hat{\mu}^*}$ and $\hat{\mu}^*$ cannot be an optimal policy for $\hat{M}$, which is a contradiction. Thus $\hat{\mu}^*$ is also an optimal policy for $M$. The proof is complete. ∎

## 5 The case with nonuniform initial exploration

The version of optimistic policy iteration described by (3) and (4) is synchronous in the sense that $n$ trajectories are simultaneously observed at each iteration, one for each initial state. It is pointed out in [12] that the scenario of picking one single state (randomly, uniformly, and independently) at each iteration to generate a trajectory from and update the cost-to-go value at this state can be captured by the following iteration:

$$J_{t+1}(i) = \begin{cases} (1 - \gamma_t)J_t(i) + \gamma_t(J^{\mu_t}(i) + w_t(i)), \\ \qquad\qquad \text{with probability } \frac{1}{n}, \\ J(i), \qquad\qquad \text{otherwise.} \end{cases} \quad (14)$$

Furthermore, this algorithm can be equivalently described in the form

$$J_{t+1} = (1 - \frac{\gamma_t}{n})J_t + \frac{\gamma_t}{n}(J^{\mu_t} + v_t), \quad (15)$$

where

$$v_t(i) = w_t(i) + (n\chi_t(i) - 1)(-J_t(i) + J^{\mu_t}(i) + w_t(i)),$$

where each $\chi_t(i)$ is a random variable satisfying $\chi_t(i) = 1$ if state $i$ is selected and $\chi_t(i) = 0$ otherwise. Then, it can be shown that $v_t$ satisfies $\mathbb{E}[v_t | \mathcal{F}_t] = 0$ and $\mathbb{E}[\|v_t\|^2 | \mathcal{F}_t] \leq A + B \|J_t\|^2$, for some constants $A$ and $B$, where we used the fact that $J^{\mu_t}$ is bounded, because there are only a finite number of policies. To use an argument similar to that in the proof of Theorem 10, one needs to show that $\mathbb{E}[\|v_t\|^2 | \mathcal{F}_t]$ is bounded.

**Proposition 14** *Under Assumptions 1 and 2, the sequence $J_t$ in (15) is bounded.*

**Proof** There exists some $D > 0$ such that $\|J^{\mu_t}\| \leq D$ for all $t \geq 0$. Boundedness of $J_t$ follows from Proposition 4.7 in [4] (see also Proposition 23 in the Appendix). ∎

Hence, $\mathbb{E}[\|v_t\|^2 | \mathcal{F}_t]$ is bounded, this time by a sequence of random variables $A_t$ that are $\mathcal{F}_t$-adapted and bounded. By a similar argument to the proof of Theorem 10, one can show that the values $J_t$ generated by (15) converge to $J^*$ under the same assumptions.

A natural question is whether we can extend (15) to the case where the states are chosen according to a nonuniform distribution such that each state has a non-zero probability of being selected. This would lead to the following update rule for $J_t$:

$$J_{t+1} = (1 - \Gamma_t)J_t + \Gamma_t(J^{\mu_t} + v_t), \quad (16)$$

where $\Gamma_t = \gamma_t \text{diag}\{p(1), p(2), \cdots, p(n)\}$ and each $p(i)$ is the probability of state $i$ being selected. It is conjectured in [12] that this may not converge (at least with the proof method therein). We also believe this is the case, although a concrete counterexample is yet to be constructed (see Section 6 for a numerical example).

Here we provide a slightly different perspective. We argue that, in the case where the states are selected nonuniformly for updating, the stepsize should be different for each state. In fact, for the classical version of MCES [11, p. 99, Chapter 5], we have the stepsize given by

$$\gamma_i(t) = \frac{1}{n_i(t)}, \quad (17)$$

where $n_i(t)$ is the number of times that state $i$ is selected up to iteration $t$, so that $J_t(i)$ is equal to the average of the simulated cumulative costs for state $i$ up to iteration $t$. This choice of stepsize is non-deterministic, but easy to implement. Alternatively, if we know *a priori* the probability of selecting each state for updating, we can design a time-varying but deterministic stepsize as

$$\gamma_t(i) = \frac{\hat{\gamma}_t}{p(i)}, \quad (18)$$

where $\hat{\gamma}_t$ is any stepsize satisfying Assumption 2. We prove that both (17) and (18) lead to convergence of the following iteration

$$J_{t+1}(i) = \begin{cases} (1 - \gamma_t(i))J_t(i) + \gamma_t(i)(J^{\mu_t}(i) + w_t(i)), \\ \qquad\qquad \text{with probability } p(i), \\ J(i), \qquad\qquad \text{otherwise,} \end{cases} \quad (19)$$

where $p(i) > 0$ is the probability that state $i$ is being selected at each iteration.

**Proposition 15** *Under Assumptions 1 and 2, both choices of stepsizes (17) and (18) lead to convergence of $J_t$ in (19) to $J^*$ in probability 1.*

**Proof** We first consider the case (18). We can equivalently write the update rule as

$$J_{t+1} = (1 - \hat{\gamma}_t)J_t + \hat{\gamma}_t(J^{\mu_t} + \hat{v}_t), \qquad (20)$$

where $\hat{\gamma}_t$ is any stepsize satisfying Assumption 2 and

$$\hat{v}_t(i) = w_t(i) + \left(\frac{\chi_t(i)}{p(i)} - 1\right)(-J_t(i) + J^{\mu_t}(i) + w_t(i)),$$

where each $\chi_t(i)$ is a random variable satisfying $\chi_t(i) = 1$ if state $i$ is selected and $\chi_t(i) = 0$ otherwise. One can easily verify that $\mathbb{E}[\hat{v}_t|\mathcal{F}_t] = 0$ and

$$\mathbb{E}[\|\hat{v}_t\|^2 \,|\mathcal{F}_t] \le \hat{A} + \hat{B} \,\|J_t\|^2,$$

where $\hat{A}$ and $\hat{B}$ are constants. Similar to Proposition 14, $J_t$ generated by (20) is bounded. The same argument as in the proof of Theorem 10 can be used to show that $J_t$ converges to $J^*$.

Now consider (17). We can write the update rule as

$$J_{t+1} = (1 - \hat{\gamma}_t)J_t + \hat{\gamma}_t(J^{\mu_t} + \hat{v}_t + \hat{u}_t), \qquad (21)$$

where $\hat{\gamma}_t = \frac{1}{t+1}$,

$$\hat{v}_t(i) = w_t(i) + \left(\frac{\chi_t(i)}{p(i)} - 1\right)\left(-J_t(i) + J^{\mu_t}(i)\right)$$
$$+ \left(\frac{(t+1)\chi_t(i)}{n_i(t)} - 1\right)w_t(i),$$

and

$$\hat{u}_t(i) = \left(\frac{t+1}{n_t(i)} - \frac{1}{p(i)}\right)\chi_t(i)(-J_t(i) + J^{\mu_t}(i)).$$

By the strong law of large numbers, $\frac{n_i(t)}{t+1} \to p(i)$ in probability 1 as $t \to \infty$. It is easy to see that there exists bounded and $\mathcal{F}_t$-adapted sequences $A_t$ and $B_t$ such that $\mathbb{E}[\hat{v}_t|\mathcal{F}_t] = 0$ and $\mathbb{E}[\|\hat{v}_t\|^2 \,|\mathcal{F}_t] \le A_t + B_t \|J_t\|^2$. Moreover, there exists an $\mathcal{F}_t$-adapted random sequence $\theta_t$ such that $\theta_t \to 0$ in probability 1, as $t \to \infty$ and $\|\hat{u}_t\| \le \theta_t(\|J_t\| + 1)$. We can then use the same argument as in the proof of Theorem 10 to show $J_t$ converges to $J^*$ in probability 1, in which we need to use Proposition 4.5 in [4] (see also Proposition 23 in the Appendix). ∎

**Remark 16** *All the convergence results obtained for the undiscounted case in this paper can be extended to the case of temporal difference $TD(\lambda)$ and model-free case (Q-learning) without much difficulty. Such results are left out due to the space limit. Interested readers should be able to refer to Sections 4 and 5 in [12] and combine the argument there with those in Sections 3–5 of this paper.*

## 6 An illustrative example

We use a simple example (adapted from [4, Example 5.11]) to illustrate the convergence behaviours of different variants of MCES.

**Example 17** *We consider a discounted problem (i.e. $\alpha < 1$) with two states and deterministic transitions show in Figure 1. Note that a discounted problem can be turned into an equivalent shortest path problem [2] by adding a terminal state and modifying the transition probability such that each transition has $1 - \alpha$ probability reaching the terminal state.*

*The states consist of $S = \{1, 2\}$ and the actions $A = \{l, r\}$. The transitions from each state under each action can be seen from Figure 1. We define the stage cost $g$ as $g(1, r) = g(2, l) = 0$ and $g(1, l) = g(2, r) = 1$. Intuitively, for each of the two states, the cost to move is 0 and the cost to stay is 1. This example was used in [4] to show possible divergence of iteration (16) when we do not restrict the frequency of selecting each of the two states for value updates.*

*There are in total four different policies $\mu_l$, $\mu_r$, $\mu_g$, and $\mu_w$, defined as follows: $\mu_l(1) = \mu_l(2) = l$, $\mu_r(1) = \mu_r(2) = r$, $\mu_g(1) = r, \mu_g(2) = l$, and $\mu_w(1) = l, \mu_w(2) = r$. The intuitive meaning of $\mu_l$ is to always move to the node on the left, while $\mu_r$ is to always move to the right. The optimal policy $\mu_g$ moves from each node to the opposite node, and the "worst" policy $\mu_w$ always stays at the current node. It is straightforward to compute the cost-to-go value for each policy as follows:*

$$J^{\mu_w} = \begin{bmatrix} \frac{1}{1-\alpha} \\ \frac{1}{1-\alpha} \end{bmatrix}, J^{\mu_g} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, J^{\mu_l} = \begin{bmatrix} \frac{1}{1-\alpha} \\ \frac{\alpha}{1-\alpha} \end{bmatrix}, J^{\mu_r} = \begin{bmatrix} \frac{\alpha}{1-\alpha} \\ \frac{1}{1-\alpha} \end{bmatrix}.$$

*Furthermore, it can be verified that $\mu_l$ is a greedy policy for a value vector $J$, if $J(1) \le J(2) - \frac{1}{\alpha}$, and $\mu_r$ is a greedy policy for a value vector $J$, if $J(2) \le J(1) - \frac{1}{\alpha}$. The optimal policy $\mu_g$ is a greedy policy for a value vector $J$, if $|J(2) - J(1)| \le \frac{1}{\alpha}$. This is also depicted in Figure 2(d), where the two black dashed lines ($J(2) = J(1) \pm \frac{1}{\alpha}$) separate the domains of the different greedy policies. Note that the optimal policy $\mu_g$ in this example does not satisfy the optimal policy feed-forward environment assumption in [13], because both states are revisited under the optimal policy.*
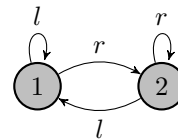


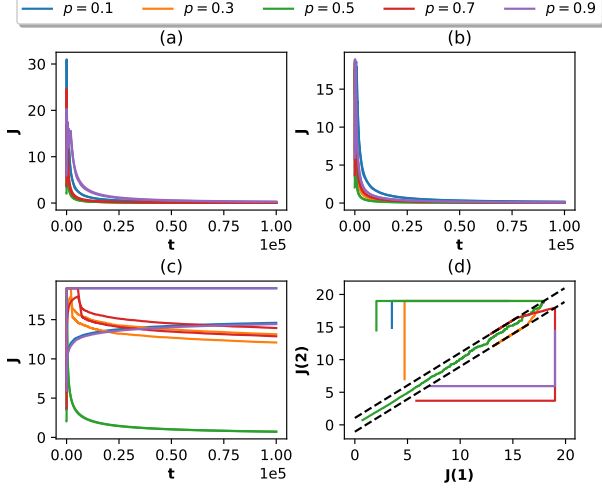Fig. 1. A two-state deterministic system.

Fig. 2. Optimistic policy iteration on Example 17 under different choices of stepsize: (a) Iteration (19 with stepsize choice (18) with $\hat{\gamma}_t = \frac{1}{t+1}$. (b) Iteration (19 with stepsize choice (17). (c)-(d) Iteration (16 with stepsize choice (18) with $\gamma_t = \frac{1}{t+1}$. The black dashed lines in (d) indicate separated domains of the three greedy policies $\mu_l$, $\mu_g$, and $\mu_r$ from top left to bottom right.

We simulate the optimistic policy iteration (19) with different probabilities $p(1) = p$ and $p(2) = 1 - p$. Figures 2(a) and 2(b) show convergence of (19) using stepsize choices (18) and (17), respectively. Results for iteration (16) are shown in Figure 2(c) and 2(d). Convergence is observed for (16) only when $p = 0.5$ (i.e., in the case of uniform selection).

## 7   Conclusions

We investigated the convergence of optimistic policy iteration, also known as Monte Carlo Exploring Starts (MCES), for the stochastic shortest path problem. These results complement known partial results on this topic and thereby help settle this long-standing open question.

There are at least two possible extensions of this work. First, the results in this paper assume that only the initial state of a simulated trajectory is picked for value updates. It would be interesting to prove convergence for the first-visit and every-visit versions of MCES [11], in which the first and every state visited on the trajectory, respectively, will be selected for value updating. Second, as pointed out in [12], it would be interesting, and perhaps very challenging, to generalize the results to situations where function approximations are used to represent values.

## Acknowledgements

## References

[1]  Robert B Ash. *Real Analysis and Probability*. Academic Press, 1972.

[2]  Dimitri P Bertsekas, Dimitri P Bertsekas, Dimitri P Bertsekas, and Dimitri P Bertsekas. *Dynamic Programming and Optimal Control*. Athena Scientific, 1995.

[3]  Dimitri P Bertsekas and John N Tsitsiklis. An analysis of stochastic shortest path problems. *Mathematics of Operations Research*, 16(3):580–595, 1991.

[4]  Dimitri P Bertsekas and John N Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, 1996.

[5]  Yuanlong Chen. On the convergence of optimistic policy iteration for stochastic shortest path problem. *arXiv preprint arXiv:1808.08763*, 2018.

[6]  Eric V Denardo. Contraction mappings in the theory underlying dynamic programming. *Siam Review*, 9(2):165–177, 1967.

[7]  Carl D Meyer. *Matrix Analysis and Applied Linear Algebra*, volume 71. SIAM, 2000.

[8]  Jacques Neveu and TP Speed. *Discrete-Parameter Martingales*, volume 10. North-Holland Amsterdam, 1975.

[9]  Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, pages 400–407, 1951.

[10]  David Silver. *Reinforcement learning and simulation-based search in computer Go*. PhD thesis, University of Alberta, 2009.

[11]  Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.

[12]  John N Tsitsiklis. On the convergence of optimistic policy iteration. *Journal of Machine Learning Research*, 3(Jul):59–72, 2002.

[13]  Che Wang and Keith Ross. On the convergence of the monte carlo exploring starts algorithm for reinforcement learning. *arXiv preprint arXiv:2002.03585*, 2020.

[14]  David Williams. *Probability with Martingales*. Cambridge University Press, 1991.

## A   Martingale Convergence Theorem

The following version of supermartingale convergence theorem stated in [4], without a proof, is widely used in convergence analysis of stochastic approximation. For completeness, we provide a self-contained proof.

**Theorem 18** *[4, Proposition 4.2] Let $\{X_t\}$, $\{Y_t\}$, and $\{Z_t\}$ be three sequences of random variables that are adapted to a filtration $\{\mathcal{F}_t\}$. Suppose that the following conditions hold:*

*(1) $X_t$, $Y_t$, and $Z_t$ are nonnegative for all $t \geq 0$.*
*(2) $\mathbb{E}[Y_{t+1} \,|\, \mathcal{F}_t] \leq Y_t - X_t + Z_t$.*
*(3) $\sum_{t=0}^{\infty} Z_t < \infty$ holds in probability 1.*

*Then $\sum_{t=0}^{\infty} X_t < \infty$ holds in probability 1 and $Y_t$ converges in probability 1 to a nonnegative random variable $Y_\infty$.*

The book [4] cited [1] and [8] for this result. However, the references [1,8] do not seem to contain an exact statement of this result, nor a proof. Here we provide a proof of this result for completeness, based on a standard version of the supermartingale convergence theorem below.

**Theorem 19** *Let $\{Y_t\}$ be a supermartingale bounded in $L^1$, i.e. $\sup_{t\geq 0} \mathbb{E}\left[|Y_t|\right] < \infty$. Then $Y_t$ converges in probability 1 to a random variable $Y_\infty$ and $Y_\infty \in L^1$.*

A proof of this result can be found, e.g., in [14, p. 109]. A variant of Theorem 19 can be proved immediately.

**Corollary 20** *Let $\{Y_t\}$ be a supermartingale. If $\sup_{t\geq 0} \mathbb{E}\left[Y_t^-\right] < \infty$, where $Y_t^-$ is the negative part of $Y_t$ defined by $Y_t = \max(0, -Y_t)$. Then $Y_t$ converges in probability 1 to a random variable $Y_\infty \in L^1$.*

**Proof** Write $|Y_t| = Y_t + 2Y_t^-$. Since $\{Y_t\}$ is a supermartingale, $\mathbb{E}[Y_t] \leq \mathbb{E}[Y_0]$ for all $t \geq 0$. Hence, $\sup_{t\geq 0} \mathbb{E}\left[Y_t^-\right] < \infty$ implies $\sup_{t\geq 0}\mathbb{E}\left[|Y_t|\right] < \infty$. The conclusion follows from Theorem 19. ■

Clearly, if $\{Y_t\}$ is a nonnegative supermartingale, then $Y_t^- \equiv 0$ and $Y_t$ converges in probability 1 according to Corollary 20.

To prove Theorem 18 based on Theorem 19, we also need the following lemma, which says that a stopped supermartingale is still a supermartingale.

**Lemma 21** *Let $\{Y_t\}$ be a supermartingale and $T$ be a stopping time. Then the stopped process $X^T := X_{T\wedge t}$, $t = 0, 1, 2, \cdots$, is still a supermartingale.*

A statement and proof of this result can be found, e.g., in [14, p. 99] or [8, p. 32].

**Proof of Theorem 18**

For each $t$, define

$$W_t = Y_t + \sum_{s=0}^{t-1} X_s - \sum_{s=0}^{t-1} Z_s.$$

It is straightforward to verify by condition (2) of Theo-

rem 18 that

$$\mathbb{E}\left[W_{t+1} \,|\, \mathcal{F}_t\right] = \mathbb{E}\left[Y_{t+1} + \sum_{s=0}^{t} X_s - \sum_{s=0}^{t} Z_s \,|\, \mathcal{F}_t\right]$$

$$= \mathbb{E}\left[Y_{t+1} \,|\, \mathcal{F}_t\right] + \sum_{s=0}^{t} X_s - \sum_{s=0}^{t} Z_s$$

$$\leq Y_t - X_t + Z_t + \sum_{s=0}^{t} X_s - \sum_{s=0}^{t} Z_s$$

$$= Y_t + \sum_{s=0}^{t-1} X_s - \sum_{s=0}^{t-1} Z_s = W_t.$$

Hence, $\{W_t\}$ is a supermartingale. For each $k \geq 0$, define a stopping time $T_k$ by

$$T_k = \inf\left\{t \geq 0 : \sum_{s=0}^{t} Z_s \geq k\right\}.$$

Then the stopped process $W_{T_k\wedge t}$, $t = 0, 1, 2, \cdots$, according to Lemma 21, is also a supermartingale. Furthermore, by the definition of $T_k$, we have $\sum_{s=0}^{T_k\wedge t-1} Z_s \leq k$, which implies $W_{T_k\wedge t} \geq -\sum_{s=0}^{T_k\wedge t-1} Z_s \geq -k$. Hence $k + W_{T_k\wedge t}$, $t = 0, 1, 2, \cdots$, is nonnegative supermartingale for each $k \geq 0$. By Corollary 20, $\lim_{t\to\infty}(k + W_{T_k\wedge t})$ exists in probability 1 for each $k \geq 0$.

Consider the event

$$\Omega_W^k = \left\{\omega \in \Omega : \lim_{t\to\infty}(k + W_{T_k\wedge t}) \text{ exists}\right\}.$$

Then $P(\Omega_W^k) = 1$ for all $k \geq 0$. Let $\Omega_W = \cap_{k=0}^{\infty}\Omega_W^k$. By continuity of probability, we have $P(\Omega_W) = 1$. Consider also the event

$$\Omega_Z = \left\{\omega \in \Omega : \sum_{t=0}^{\infty} Z_t(\omega) < \infty\right\}.$$

Then $P(\Omega_Z) = 1$. It follows that $P(\Omega_Z \cap \Omega_W) = 1$.

Consider any $\omega \in \Omega_Z \cap \Omega_W$. Since $\omega \in \Omega_Z$, there exists some $k \geq 0$ such that $\sum_{t=0}^{\infty} Z_t(\omega) < k$. Hence, for this $k$, we have $T_k(\omega) = \infty$. Since $\omega \in \Omega_W$, $\lim_{t\to\infty} W_t(\omega) = \lim_{t\to\infty}(k + W_{T_k(\omega)\wedge t}((\omega)))$ exists. We have proved that $W_t$ converges in probability 1. By the definition of $W_t$ and the fact that $\sum_{t=0}^{\infty} Z_t < \infty$ in probability 1, $Y_t + \sum_{s=0}^{t-1} X_s$ converges in probability 1. Since $X_t$ is nonnegative, condition (2) of Theorem also holds with $X_t \equiv 0$. Hence, repeating the argument above with $X_t \equiv 0$ would show that $Y_t$ converges in probability 1. This in turn implies $\sum_{t=0}^{\infty} X_t < \infty$ in probability 1. ■

## B  Stochastic Approximation

Based on the supermartingale convergence theorem, in this section, we provide a more straightforward proof of the convergence result on stochastic approximation arguments we used in this paper.

Consider a sequence $\{J_t\}$ generated using the update rule

$$J_{t+1}(i) = (1 - \gamma_t(i))J_t(i) + \gamma_t(H_t J_t(i) + w_t(i) + u_t(i)), \tag{B.1}$$

where $\gamma_t$, $H_t$, $w_t$, and $u_t$ satisfy the following.

**Assumption 22**  *We have*

(1) *$\sum_{t=0}^{\infty} \gamma(i) = \infty$ and $\sum_{t=0}^{\infty} \gamma^2(i) < \infty$ for all $i$.*
(2) *There exists a positive vector $\theta \in \mathbb{R}^n$, a vector $J^* \in \mathbb{R}^n$, and sclars $\beta \in [0,1)$ and $D \geq 0$ such that $\|H_t J - J^*\|_\theta \leq \beta \|J - J^*\|_\theta + D$. We also assume that $H_t J_t$ is $\mathcal{F}_t$-adapted.*
(3) *There exist constants $A$ and $B$ such that*

$$\mathbb{E}\left[w_t(i)\,|\,\mathcal{F}_t\right] = 0, \quad \mathbb{E}\left[w_t^2(i)\,|\,\mathcal{F}_t\right] \leq A_t + B_t \|J_t\|^2$$

*for all $i$ and $t$, where $\|\cdot\|$ is any norm and $A_t$ and $B_t$ are $\mathcal{F}_t$-adapted and bounded.*
(4) *There exists an $\mathcal{F}_t$-adapted random sequence $\theta_t$ such that $\lim_{t \to \infty} \theta_t = 0$ in probability 1 and $|u_t(i)| \leq \theta_t(\|J_t\| + 1)$ for all $i$ and $t$, where $\|\cdot\|$ is any norm. We also assume that $u_t$ is $\mathcal{F}_t$-adapted.*

The following result is essentially Propositions 4.7 and 4.5 in [4] combined together. Here we provide a more direct proof from the supermartingale convergence theorem.

**Proposition 23**  *Let $J_t$ be generated by (B.1). Suppose that Assumption 22 holds. Then*

(1) *$J_t$ is bounded in probability 1, and*
(2) *$J_t$ converges to $J^*$ in probability 1 if $D = 0$.*

Since the analysis with the weighted maximum norm $\|\cdot\|_\theta$ is very similar to that of the maximum $\|\cdot\|_\infty$. In the following proofs, we only consider the maximum norm and denote it by $\|\cdot\|$.

**Lemma 24**  *Consider*

$$V_{t+1}(i) = (1 - \gamma_t(i))V_t(i) + \gamma_t(i)w_t(i), \quad t \geq t_0 \geq 0.$$

(1) *If (3) of Assumption 22 holds with $B = 0$ (or $J_t$ is bounded), then $V_t$ converges to 0 in probability 1.*
(2) *Let $G_t$ be a nondecreasing $\mathcal{F}_t$-adapted scalar random sequence such that $G_t \geq \mu \|J_t\| + \nu$ for all $t \geq t_0$, where $\mu$ and $\nu$ are positive constants. Then $\frac{V_t}{G_t}$ converges to 0 in probability 1.*

**Proof**  We prove (2) first. The proof for (1) is a special case. We have

$$V_{t+1}^2(i) = (1 - \gamma_t(i))^2 V_t^2(i) + 2\gamma_t(i)(1 - \gamma_t(i))V_t(i)w_t(i) + \gamma_t^2(i)w_t^2(i).$$

Since $G_t$ is a nondecreasing sequence, we obtain

$$\frac{V_{t+1}^2(i)}{G_{t+1}^2} \leq \frac{V_{t+1}^2(i)}{G_t^2}$$
$$= (1 - \gamma_t(i))^2 \frac{V_t^2(i)}{G_t^2} + 2\gamma_t(i)(1 - \gamma_t(i))\frac{V_t(i)w_t(i)}{G_t^2}$$
$$+ \gamma_t^2(i)\frac{w_t^2(i)}{G_t^2}.$$

Taking condition expectation from both sides and noticing that $\mathbb{E}\left[w_t(i)\,|\,\mathcal{F}_t\right] = 0$ and $V_t$ and $G_t$ are adapted to $\mathcal{F}_t$ and independent of $w_t$, we have

$$\mathbb{E}\left[\frac{V_{t+1}^2(i)}{G_{t+1}^2}\,\bigg|\,\mathcal{F}_t\right] \leq (1 - \gamma_t(i))^2 \frac{V_t^2(i)}{G_t^2} + \gamma_t^2(i)\frac{\mathbb{E}\left[w_t^2(i)\right]}{G_t^2}$$
$$\leq (1 - \gamma_t(i))^2 \frac{V_t^2(i)}{G_t^2} + \gamma_t^2(i)\frac{A_t + B_t \|J_t\|^2}{G_t^2}$$
$$\leq (1 - 2\gamma_t(i) + \gamma_t^2(i))\frac{V_t^2(i)}{G_t^2} + \gamma_t^2(i)K_t,$$

for some $\mathcal{F}_t$-adapted and bounded $K_t$, where we used (3) of Assumption 22 and the fact that $G_t \geq \mu \|J_t\| + \nu$ for all $t$.

Since $\gamma_t(i) \to 0$ as $t \to \infty$, for $t$ sufficiently large, we have $\gamma_t^2(i) \leq \gamma_t(i)$ and

$$\mathbb{E}\left[\frac{V_{t+1}^2(i)}{G_{t+1}^2}\,\bigg|\,\mathcal{F}_t\right] \leq \frac{V_t^2(i)}{G_t^2} - \gamma_t(i)\frac{V_t^2(i)}{G_t^2} + \gamma_t^2(i)K_t.$$

Let $Y_t = \frac{V_t^2(i)}{G_t^2}$, $X_t = \gamma_t(i)\frac{V_t^2(i)}{G_t^2}$, and $Z_t = \gamma_t^2(i)K_t$. Since $K_t$ is bounded, we have $\sum_t^\infty \gamma_t^2(i)K_t < \infty$ in probability 1. Then the conditions of Theorem 18 are satisfied for $t$ sufficiently large. By Theorem 18, $Y_t = \frac{V_t^2(i)}{G_t^2}$ converges in probability 1 and $\sum_{t=0}^\infty \gamma_t(i)\frac{V_t^2(i)}{G_t^2} < \infty$ in probability 1, which in turn implies $Y_t = \frac{V_t^2(i)}{G_t^2}$ converges in probability 1 to 0, because otherwise we would have $\sum_{t=0}^\infty \gamma_t(i)\frac{V_t^2(i)}{G_t^2} = \infty$ since $\sum_{t=0}^\infty \gamma_t(i) = \infty$. To prove (1), note that if $B = 0$ (or $J_t$ is bounded), we can set $G_t = 1$ and prove convergence of $V_t$ in the same way. ∎

The first part of the above lemma is Corollary 4.1 in [4], for which we provide a more direct proof here. The second part appears to be new.

**Lemma 25** *Consider*

$$Y_{t+1}(i) = (1 - \gamma_t(i))Y_t(i) + \gamma_t(i)G_t, \quad t \geq t_0,$$

*where $t_0 \geq 0$ and $G_t$ is a positive nondecreasing scalar random sequence. Then $\limsup_{t \to \infty} \left| \frac{Y_t(i)}{G_t} \right| \leq 1$.*

**Proof** We have

$$\left| \frac{Y_{t+1}(i)}{G_{t+1}} \right| \leq |1 - \gamma_t(i)| \left| \frac{Y_t(i)}{G_{t+1}} \right| + \gamma_t(i)$$

$$\leq (1 - \gamma_t(i)) \left| \frac{Y_t(i)}{G_t} \right| + \gamma_t(i),$$

for $t$ sufficiently large such that $\gamma_t(i) < 1$. Consider the iteration
$$Z_{t+1} = (1 - \gamma_t(i))Z_t + \gamma_t(i),$$
with $Z_{t_0} = \frac{Y_{t_0}}{G_{t_0}}$. It is easy to verify that $Z_t \to 1$, as $t \to \infty$ (this can in fact be seen as a special case of Lemma 24 with $V_t = Z_t - 1$ and $w_t = 0$). By comparison, $\frac{Y_t}{G_t} \leq Z_t$ for all $t \geq t_0$. Hence, $\limsup_{t \to \infty} \left| \frac{Y_t(i)}{G_t} \right| \leq 1$. ∎

**Proof of Proposition 23**

Note that all the estimates on random variables in this proof are meant to hold in probability 1.

Fix an $\eta \in (0, 1)$ such that $\beta + 2\eta < 1$. Since $\theta_t \to 0$, as $t \to \infty$, for any $\varepsilon > 0$, there exists $t_0 = t_0(\varepsilon)$ such that $\theta_t < \varepsilon$ for all $t \geq t_0$. Since $\gamma_t(i) \to 0$, we can also assume $t_0$ is picked sufficiently large such that $\gamma_t(i) < 1$.

By Assumption 22, we have

$$\|H_t J_t\| + \theta_t(\|J_t\| + 1) \leq \beta \|J_t\| + D + \varepsilon(\|J_t\| + 1)$$
$$\leq (\beta + \varepsilon) \sup_{0 \leq s \leq t} \|J_s\| + (D + \varepsilon)$$
$$= G_t, \quad \forall t \geq t_0, \qquad (B.2)$$

where $G_t := (\beta + \varepsilon) \sup_{0 \leq s \leq t} \|J_s\| + (D + \varepsilon)$. Then $G_t$ satisfies the assumptions in Lemmas 24 and 25. Now consider

$$Y_{t+1}(i) = (1 - \gamma_t(i))Y_t(i) + \gamma_t(i)G_t, \quad t \geq t_0,$$

and

$$V_{t+1}(i) = (1 - \gamma_t(i))V_t(i) + \gamma_t(i)w_t(i), \quad t \geq t_0.$$

Set $Y_{t_0} = J_{t_0}$ and $V_{t_0} = 0$.

**Claim:** We have $Y_t - V_t \leq J_t \leq Y_t + V_t$ for all $t \geq t_0$.

**Proof of the claim:** We prove it by induction. For $t = t_0$, we have $J_{t_0} = Y_{t_0} + V_{t_0}$. Assume the inequality holds

for some $t \geq t_0$. By (B.2), we have

$$J_{t+1}(i) \leq (1 - \gamma_t(i))J_t(i) + \gamma_t(i)(H_t J_t(i) + w_t(i) + u_t(i))$$
$$\leq (1 - \gamma_t(i))(Y_t(i) + V_t(i)) + \gamma_t(i)G_t + \gamma_t(i)w_t(i)$$
$$= Y_{t+1}(i) + V_{t+1}(i).$$

The other half of the inequality similarly holds. ∎

By Lemma 24, we have $\frac{V_t(i)}{G_t} \to 0$, as $t \to \infty$. By Lemma 25, $\limsup_{t \to \infty} \left| \frac{Y_t(i)}{G_t} \right| \leq 1$. For the same $\varepsilon > 0$, there exists $T \geq t_0$ such that $|V_t(i)| \leq \varepsilon G_t$ and $|Y_t(i)| \leq (1 + \varepsilon)G_t$ for all $t \geq T$. Hence the above claim implies that $\|J_t\| \leq (1 + 2\varepsilon)G_t$ for all $t \geq T$. In view of the definition $G_t$ from (B.2), we obtain

$$\|J_t\| \leq (1 + 2\varepsilon)[(\beta + \varepsilon) \sup_{0 \leq s \leq t} \|J_s\| + (D + \varepsilon)], \quad (B.3)$$

for all $t \geq T$. Fix $\varepsilon \in (0, 1)$ sufficiently small such that $\mu := (1 + 2\varepsilon)(\beta + \varepsilon) < 1$. It follows that

$$\sup_{0 \leq s \leq t} \|J_s\| \leq \mu \sup_{0 \leq s \leq t} \|J_s\| + C, \quad \forall t \geq 0,$$

where $C = \max((1 + 2\varepsilon)(D + \varepsilon), \sup_{0 \leq s \leq T} \|J_s\|)$. Hence, we obtain an explicit bound for $J_t$ as

$$\|J_t\| \leq \sup_{0 \leq s \leq t} \|J_s\| \leq \frac{C}{1 - \mu}, \quad \forall t \geq 0. \qquad (B.4)$$

Note that $C$ is a random variable. This proves item (1).

We now prove item (2). Let $C_0 = \frac{C}{1 - \mu}$ and $T_0 = T$. Then $C_0$ is $\mathcal{F}_T$-measurable. For any $\varepsilon_0 \in (0, \varepsilon)$, define $G_t = (\beta + \varepsilon)C_0 + \varepsilon$ for $t \geq T_0$. Then (B.2) holds with this $G_t$. By repeating the argument preceding (B.3), we can show that there exists some $T_1 \geq T_0$ such that

$$\|J_t\| \leq (1 + 2\varepsilon_0)[(\beta + \varepsilon_0)C_0 + \varepsilon_0)], \quad \forall t \geq T_1. \quad (B.5)$$

We can pick $\varepsilon_0$ sufficiently small such that

$$\|J_t\| \leq (\beta + \eta)C_0, \quad \forall t \geq T_1.$$

We can inductively show that there exists a sequence $\{T_k\}$ such that

$$\|J_t\| \leq (\beta + \eta)^k C_0, \quad \forall t \geq T_k.$$

Hence $J_t \to 0$ as $t \to \infty$. This proves item (2). ∎