



ON A VERIFICATION THEOREM FOR OPTIMALITY OF STRATEGIES IN MARKOV DECISION PROCESSES WITH TOTAL COST

ALEXEY PIUNOVSKIY^{1,*}, YI ZHANG²

¹Department of Mathematical Sciences, University of Liverpool, Liverpool, UK

²School of Mathematics, University of Birmingham, Edgbaston, Birmingham, B15 2TT, UK

Dedicated to the memory of Professor Uriel Rothblum

Abstract. The goal of this brief article is to present a verification theorem, which can be used to show the optimality of a given strategy, without computing the Bellman function in advance, for a Markov decision process (MDP) with a general-signed cost function. An example is elaborated on, which shows that one of the minor conditions in the verification theorem is really important.

Keywords. Markov decision process; Total cost; Optimal strategy.

2020 Mathematics Subject Classification. 90C40, 49L20.

1. INTRODUCTION

Markov decision processes (MDPs) is a well developed branch of applied probability, statistics and operational research with numerous applications to real-life problems. The current paper is devoted to the MDPs with a total expected cost. One can say that there are two main approaches, in some sense dual: convex analytic approach and dynamic programming (DP). As for the first one, a substantial progress was made in the recent years by Professor U. Rothblum [4].

Professor U. Rothblum was also active in developing the dynamic programming approach; see, e.g., [5]. One can find many useful statements on DP in the famous monographs by Bertsekas and Shreve [1] and by Puterman [10]. Usually, one of the main steps is calculating the Bellman function, given by the minimal possible expected cost with any given initial state. After that, a control strategy is optimal if it is so called conserving and equalizing; see, e.g. [1, Prop. 9.12] or [10, Thm. 7.1.7]. Both conditions are stated in terms of the Bellman function. There are fewer theorems which help to check the optimality of a strategy without calculating the Bellman function *a priori*. For example, Proposition 9.13 of [1] is about only negative models and discounted ones with the bounded cost function; [6, Thm. 4.2.3(c)] is about discounted MDPs with an unbounded cost function. The goal of the current paper is to prove a similar statement in the more general case: no restrictions on the sign of the cost function, and the MDP is not necessarily absorbing.

*Corresponding author.

E-mail address: piunov@liverpool.ac.uk (A. Piunovskiy), y.zhang.29@bham.ac.uk (Y. Zhang).

Received November 9, 2023; Accepted March 14, 2024.

In Sections 2 and 3, we describe the MDP under investigation, formulate and prove our main result, Theorem 3.1. In Section 4, we elaborate on a counter-example showing that one of the conditions in that theorem cannot be omitted. All over the text, when calculating expectations, we accept the following conventions:

$$E[Y] := E[Y^+] + E[Y^-], \text{ where } \infty - \infty := \infty \text{ and } Y^+ := \max\{Y, 0\}, Y^- := \min\{Y, 0\}.$$

2. DESCRIPTION OF THE MODEL

The primitives of an MDP are the following.

- The state and action spaces \mathbf{X} and \mathbf{A} are nonempty Borel spaces, endowed with their σ -algebras $\mathcal{B}(\mathbf{X})$ and $\mathcal{B}(\mathbf{A})$.
- The transition probability $p(dy|x, a)$ is a stochastic kernel from $\mathbf{X} \times \mathbf{A}$ to $\mathcal{B}(\mathbf{X})$.
- The $[-\infty, +\infty]$ -valued one-step cost function $c(\cdot, \cdot)$ on $\mathbf{X} \times \mathbf{A}$.

The initial distribution P_0 on $\mathcal{B}(\mathbf{X})$ is fixed. $\Omega := (\mathbf{X} \times \mathbf{A})^\infty$ is the sample space, endowed with the product σ -algebra \mathcal{F} . The component-wise projections

$$\Omega \ni (x_0, a_1, x_1, \dots) \rightarrow x_t, t = 0, 1, 2, \dots$$

and

$$\Omega \ni (x_0, a_1, x_1, \dots) \rightarrow a_t, t = 1, 2, \dots$$

are denoted by X_t and A_t , and $\{X_t\}_{t=0}^\infty$, $\{A_t\}_{t=1}^\infty$ are the controlled and the controlling random processes.

Definition 2.1 (Strategy). A control strategy $\pi = \{\pi_n\}_{n=1}^\infty$ is a sequence of stochastic kernels such that for each $n = 1, 2, \dots$, $\pi_n(da|x_0, a_1, \dots, x_{n-1})$ is a stochastic kernel from $(\mathbf{X} \times \mathbf{A})^{n-1} \times \mathbf{X}$ to $\mathcal{B}(\mathbf{A})$, where $(\mathbf{X} \times \mathbf{A})^0 \times \mathbf{X} := \mathbf{X}$. A strategy is called stationary if, for all $n = 1, 2, \dots$, $\pi_n(da|x_0, a_1, \dots, x_{n-1}) = \pi(da|x_{n-1})$ is the same stochastic kernel from \mathbf{X} to $\mathcal{B}(\mathbf{A})$. If for a stationary strategy π , there is a measurable mapping φ from \mathbf{X} to \mathbf{A} such that $\pi(da|x) = \delta_{\varphi(x)}(da)$ for all $x \in \mathbf{X}$, where $\delta_{\varphi(x)}(da)$ is the Dirac measure concentrated on the singleton $\{\varphi(x)\}$, then π is called a deterministic stationary strategy, identified with the (measurable) selector φ .

As is well known, for each control strategy π and initial distribution P_0 , there is a unique strategic measure on the sample space (Ω, \mathcal{F}) , denoted as $P_{P_0}^\pi$, which is specified by the following conditions:

$$P_{P_0}^\pi(X_0 \in dy) = P_0(dy);$$

and for each $n = 1, 2, \dots$, $\Gamma_i^{\mathbf{X}} \in \mathcal{B}(\mathbf{X})$ ($i = 0, 1, \dots, n$) and $\Gamma_i^{\mathbf{A}} \in \mathcal{B}(\mathbf{A})$ ($i = 1, 2, \dots, n$),

$$\begin{aligned} & P_{P_0}^\pi(X_0 \in \Gamma_0^{\mathbf{X}}, A_1 \in \Gamma_1^{\mathbf{A}}, \dots, X_{n-1} \in \Gamma_{n-1}^{\mathbf{X}}, A_n \in \Gamma_n^{\mathbf{A}}) \\ &= \int_{\Gamma_0^{\mathbf{X}} \times \Gamma_1^{\mathbf{A}} \times \dots \times \Gamma_{n-1}^{\mathbf{X}}} \pi_n(\Gamma_n^{\mathbf{A}}|x_0, a_1, \dots, x_{n-1}) P_{P_0}^\pi(X_0 \in dx_0, A_1 \in da_1, \dots, X_{n-1} \in dx_{n-1}); \end{aligned}$$

and

$$\begin{aligned} & P_{P_0}^\pi(X_0 \in \Gamma_0^{\mathbf{X}}, A_1 \in \Gamma_1^{\mathbf{A}}, \dots, X_n \in \Gamma_n^{\mathbf{X}}) \\ &= \int_{\Gamma_0^{\mathbf{X}} \times \Gamma_1^{\mathbf{A}} \times \dots \times \Gamma_{n-1}^{\mathbf{X}} \times \Gamma_n^{\mathbf{A}}} p(\Gamma_n^{\mathbf{X}}|x_{n-1}, a_n) \\ & \quad \times P_{P_0}^\pi(X_0 \in dx_0, A_1 \in da_1, \dots, X_{n-1} \in dx_{n-1}, A_n \in da_n). \end{aligned}$$

For details, see [3, 6, 8]. The expectation taken with respect to $P_{P_0}^\sigma$ is denoted by $E_{P_0}^\sigma$. In case $P_0(dx) = \delta_{x_0}(dx)$ is a Dirac measure, we use notations $P_{x_0}^\pi$ and $E_{x_0}^\pi$.

The problem is about minimizing the total expected cost:

$$v_{P_0}^\pi := \limsup_{T \rightarrow \infty} \mathbb{E}_{P_0}^\pi \left[\sum_{t=1}^T c(X_{t-1}, A_t) \right] \rightarrow \inf_{\pi} \quad (2.1)$$

Here and below,

$$\mathbb{E}_{P_0}^\pi \left[\sum_{t=1}^T c(X_{t-1}, A_t) \right] := \mathbb{E}_{P_0}^\pi \left[\sum_{t=1}^T c^+(X_{t-1}, A_t) \right] + \mathbb{E}_{P_0}^\pi \left[\sum_{t=1}^T c^-(X_{t-1}, A_t) \right]. \quad (2.2)$$

Recall that convention $\infty - \infty := \infty$ is in use here.

Usually, it is assumed that the expression $\infty - \infty$ never appears in (2.2), or, more specifically, when considering only Dirac measures P_0 , it is assumed that

$$\inf_{\pi} \mathbb{E}_{x_0}^\pi \left[\sum_{t=1}^{\infty} c^-(X_{t-1}, A_t) \right] > -\infty, \quad \forall x_0 \in \mathbf{X}, \quad (2.3)$$

see [7, Ass. 9.3.2]. Under this and other mild conditions, a very strong theorem characterizing the optimal strategies was proved in [7]; see Theorem 9.5.5 therein. This and similar statements in [1, 10] assume that the Bellman function

$$v_x^* := \inf_{\pi} v_x^\pi, \quad x \in \mathbf{X}$$

was preliminary calculated. In the next section, we formulate and prove our main result which allows to check the optimality of a strategy without calculating function v_x^* . We also compare our result with the similar previously published ones.

3. MAIN RESULT

Theorem 3.1. *Suppose that $\hat{\pi}$ is a strategy such that $v_{P_0}^{\hat{\pi}}$ and $v_x^{\hat{\pi}}$ for all $x \in \mathbf{X}$ are finite, and equation*

$$v_x^{\hat{\pi}} = \inf_{a \in \mathbf{A}} \left\{ c(x, a) + \int_{\mathbf{X}} v_y^{\hat{\pi}} p(dy|x, a) \right\}, \quad x \in \mathbf{X}. \quad (3.1)$$

is satisfied. Assume that, for each strategy π , $\mathbb{E}_{P_0}^\pi[v_{X_t}^{\hat{\pi}}]$ is finite for all $t = 0, 1, \dots$ and

$$\liminf_{T \rightarrow \infty} \mathbb{E}_{P_0}^\pi[v_{X_T}^{\hat{\pi}}] \leq 0. \quad (3.2)$$

Then $\hat{\pi}$ is optimal for the initial distribution P_0 .

Proof. The target is to show that, for any strategy π , $v_{P_0}^\pi \geq v_{P_0}^{\hat{\pi}}$. This inequality is obvious in the case $v_{P_0}^\pi = +\infty$, so below we assume that $v_{P_0}^\pi < \infty$.

Since

$$\begin{aligned} \mathbb{E}_{P_0}^\pi \left[\int_{\mathbf{X}} (v_y^{\hat{\pi}})^- p(dy|X_t, A_{t+1}) \right] &= \mathbb{E}_{P_0}^\pi [(v_{X_{t+1}}^{\hat{\pi}})^-] > -\infty, \quad t = 0, 2, \dots, \\ \mathbb{E}_{P_0}^\pi \left[\int_{\mathbf{X}} v_y^{\hat{\pi}} p(dy|X_t, A_{t+1}) \right] &= \mathbb{E}_{P_0}^\pi [v_{X_{t+1}}^{\hat{\pi}}] \in \mathbb{R} \end{aligned}$$

for each $t = 0, 1, \dots, T-1$.

Now equation (3.1) implies that, for each $T = 1, 2, \dots$,

$$\mathbb{E}_{P_0}^\pi \left[\sum_{t=1}^T \left\{ c(X_{t-1}, A_t) - v_{X_{t-1}}^{\hat{\pi}} + \int_{\mathbf{X}} v_y^{\hat{\pi}} p(dy|X_{t-1}, A_t) \right\} \right] \geq 0. \quad (3.3)$$

The left-hand side of (3.3) takes the form of the sum of expectations:

$$\begin{aligned}
& \mathbb{E}_{\mathbb{P}_0}^\pi \left[\sum_{t=1}^T c(X_{t-1}, A_t) \right] - \mathbb{E}_{\mathbb{P}_0}^\pi [v_{X_0}^{\hat{\pi}}] + \mathbb{E}_{\mathbb{P}_0}^\pi \left[\int_{\mathbf{X}} v_y^{\hat{\pi}} p(dy|X_0, A_1) \right] \\
& - \mathbb{E}_{\mathbb{P}_0}^\pi [v_{X_1}^{\hat{\pi}}] + \mathbb{E}_{\mathbb{P}_0}^\pi \left[\int_{\mathbf{X}} v_y^{\hat{\pi}} p(dy|X_1, A_2) \right] - \dots + \mathbb{E}_{\mathbb{P}_0}^\pi \left[\int_{\mathbf{X}} v_y^{\hat{\pi}} p(dy|X_{T-1}, A_T) \right] \\
& = \mathbb{E}_{\mathbb{P}_0}^\pi \left[\sum_{t=1}^T c(X_{t-1}, A_t) \right] - v_{\mathbb{P}_0}^{\hat{\pi}} + \mathbb{E}_{\mathbb{P}_0}^\pi [v_{X_T}^{\hat{\pi}}] \geq 0.
\end{aligned} \tag{3.4}$$

For each $T \geq 1$, $\mathbb{E}_{\mathbb{P}_0}^\pi \left[\sum_{t=1}^T c(X_{t-1}, A_t) \right]$ cannot equal $-\infty$ in view of the last inequality, where $v_{\mathbb{P}_0}^{\hat{\pi}}$ and $\mathbb{E}_{\mathbb{P}_0}^\pi [v_{X_T}^{\hat{\pi}}]$ are finite. Hence, in view of (2.2), for each $T \geq 1$, $\mathbb{E}_{\mathbb{P}_0}^\pi \left[\sum_{t=1}^T c^-(X_{t-1}, A_t) \right] > -\infty$. Also, for each $T \geq 1$, $\mathbb{E}_{\mathbb{P}_0}^\pi \left[\sum_{t=1}^T c(X_{t-1}, A_t) \right]$ is different from ∞ because we assumed that $v_{\mathbb{P}_0}^\pi < \infty$. In greater detail, if, for some $N \geq 1$, $\mathbb{E}_{\mathbb{P}_0}^\pi \left[\sum_{t=1}^N c(X_{t-1}, A_t) \right] = \infty$, then, in view of (2.2), $\mathbb{E}_{\mathbb{P}_0}^\pi \left[\sum_{t=1}^N c^+(X_{t-1}, A_t) \right] = \infty$, so that for all $T \geq N$, $\mathbb{E}_{\mathbb{P}_0}^\pi \left[\sum_{t=1}^T c^+(X_{t-1}, A_t) \right] = \infty$ and, again by (2.2) and the previous observation, $\mathbb{E}_{\mathbb{P}_0}^\pi \left[\sum_{t=1}^T c(X_{t-1}, A_t) \right] = \infty$. This would contradict the assumption of $v_{\mathbb{P}_0}^\pi < \infty$. Thus, all terms in (3.4) are finite. After passing to the limit, we conclude that

$$\begin{aligned}
v_{\mathbb{P}_0}^\pi &= \limsup_{T \rightarrow \infty} \mathbb{E}_{\mathbb{P}_0}^\pi \left[\sum_{t=1}^T c(X_{t-1}, A_t) \right] \geq v_{\mathbb{P}_0}^{\hat{\pi}} + \limsup_{T \rightarrow \infty} \{-\mathbb{E}_{\mathbb{P}_0}^\pi [v_{X_T}^{\hat{\pi}}]\} \\
&= v_{\mathbb{P}_0}^{\hat{\pi}} - \liminf_{T \rightarrow \infty} \mathbb{E}_{\mathbb{P}_0}^\pi [v_{X_T}^{\hat{\pi}}] \geq v_{\mathbb{P}_0}^{\hat{\pi}}.
\end{aligned}$$

□

Before we compare Theorem 3.1 with other existing results and provide a corollary, let us recall several definitions.

- Definition 3.2.** (a) An MDP is called absorbing at a specific (cemetery) state $\Delta \in \mathbf{X}$ under the initial distribution \mathbb{P}_0 if $p(\{\Delta\}|\Delta, a) \equiv 1$, $c(\Delta, a) \equiv 0$ and $\mathbb{E}_{\mathbb{P}_0}^\pi [T_0] < \infty$ for all strategies π , where $T_0 := \inf\{t \geq 0 : X_t = \Delta\}$ is the time to absorption, having accepted that the infimum over the empty set equals ∞ .
- (b) An MDP, which is absorbing at Δ under each initial distribution \mathbb{P}_0 , is called discounted if $p(\{\Delta\}|x, a) \equiv 1 - \beta$ for all $x \neq \Delta$, where $\beta \in (0, 1)$ is the ‘‘discount factor’’. (Note that the requirement $\mathbb{E}_{\mathbb{P}_0}^\pi [T_0] < \infty$ is satisfied automatically for all \mathbb{P}_0 as soon as $p(\{\Delta\}|x, a) \equiv 1 - \beta$ with $\beta \in (0, 1)$.)

If the model is negative (i.e., $c(\cdot, \cdot) \leq 0$) or discounted with a bounded cost function $c(\cdot, \cdot)$, then the following version of Theorem 3.1 is known; see e.g., [1, Prop. 9.13]: a stationary strategy $\hat{\pi}$ is optimal (for all initial states $x \in \mathbf{X}$) if and only if equation (3.1) holds.

For the discounted model with an unbounded cost function $c(\cdot, \cdot)$, under some additional conditions (e.g., $c(x, \cdot) \geq 0$ is inf-compact on \mathbf{A} for each $x \in \mathbf{X}$, etc) the following version of Theorem 3.1 is known in e.g., [6, Thm. 4.2.3(c)]: if $\lim_{T \rightarrow \infty} \mathbb{E}_x^\pi [v_{X_T}^{\hat{\pi}}] = 0$ for all $x \in \mathbf{X}$ and all strategies π such that $v_x^\pi < \infty$, then $\hat{\pi}$ is optimal (for all initial states $x \in \mathbf{X}$) if and only if equation (3.1) holds.

In the general model, under some additional conditions (e.g., inequality (2.3) etc), if $\limsup_{T \rightarrow \infty} \mathbb{E}_x^\pi [v_{X_T}^{\hat{\pi}}] = 0$ for all $x \in \mathbf{X}$ and all strategies π , then $\hat{\pi}$ is optimal (for all initial states $x \in \mathbf{X}$) if and only if equation (3.1) holds: see [7, Thm. 9.5.13].

Corollary 3.3 (from Theorem 3.1). *Suppose $\sup_{(x,a) \in \mathbf{X} \times \mathbf{A}} |c(x, a)| \leq K < \infty$ and the model is absorbing at Δ under each initial state $x \in \mathbf{X}$.*

- (a) If the model is absorbing at Δ also under the initial distribution P_0 and equations (3.1) and (3.2) are satisfied, then the strategy $\hat{\pi}$ is optimal for P_0 .
- (b) If equation (3.1) is satisfied and (3.2) holds for $P_0(dx) = \delta_{x_0}(dx)$ for all $x_0 \in \mathbf{X}$, then the strategy $\hat{\pi}$ is optimal for all initial states $x \in \mathbf{X}$.

Remark 3.4. The requirement that the model is absorbing at Δ also under the initial distribution P_0 is independent, because there are examples when absorbing under each initial state $x \in \mathbf{X}$ model is not absorbing under another initial distribution P_0 . See, e.g., Example 2.2.2 or 2.2.21 in [9].

Proof of Corollary 3.3. (a) It is known that, for the absorbing model, $E_{P_0}^\pi [T_0] \leq M < \infty$ for all strategies π , where the constant M can depend on P_0 : see [4, p.132]. Therefore, for all strategies π , $\left| E_{P_0}^\pi [\sum_{t=1}^\infty c^\pm(X_{t-1}, A_t)] \right| \leq KM$ and $E_{P_0}^\pi [\sum_{t=1}^\infty c(X_{t-1}, A_t)]$ is finite and coincides with $v_{P_0}^\pi$ (see (2.1)) for all strategies π .

The same reasoning holds for $P_0(dx) = \delta_{x_0}(dx)$, for each $x_0 \in \mathbf{X}$. Thus, $v_{P_0}^{\hat{\pi}}$ and $v_x^{\hat{\pi}}$ for all $x \in \mathbf{X}$ are finite.

Now, for each strategy π ,

$$E_{P_0}^\pi [v_{X_t}^{\hat{\pi}}] = E_{P_0}^\pi \left[E_{X_t}^{\hat{\pi}} \left[\sum_{n=1}^\infty c(X_{n-1}, A_n) \right] \right] = E_{P_0}^{\tilde{\pi}} \left[\sum_{\tau=t+1}^\infty c(X_{\tau-1}, A_\tau) \right],$$

where $\tilde{\pi}$ is the following combination of the strategies π and $\hat{\pi}$:

$$\tilde{\pi}_n(da|x_0, a_1, \dots, x_{n-1}) = \begin{cases} \pi_n(da|x_0, a_1, \dots, x_{n-1}), & \text{if } n \leq t; \\ \hat{\pi}_{n-t}(da|x_t, a_{t+1}, \dots, x_{n-1}), & \text{if } n > t. \end{cases}$$

Thus, again $\left| E_{P_0}^\pi [v_{X_t}^{\hat{\pi}}] \right| \leq KM$ is finite, and assertion (a) follows from Theorem 3.1.

(b) After we put $P_0(dx) = \delta_{x_0}(dx)$ for an arbitrarily fixed $x_0 \in \mathbf{X}$, the strategy $\hat{\pi}$ is optimal for the initial state x_0 according to (a). \square

4. COUNTER-EXAMPLE

Example 3.1.4 of [2] presents a non-optimal strategy $\hat{\pi}$ in a discounted positive model (when $c(\cdot, \cdot) \geq 0$) satisfying equation (3.1), for which $v_x^{\hat{\pi}} = \infty$ and $\lim_{T \rightarrow \infty} E_x^\pi [v_{X_T}^{\hat{\pi}}] = \infty$ for some states x and some strategies π . Below, we elaborate on another illustrative example, based on [9, Subsec. 3.2.3], where all functions are finite, which confirms that requirement (3.2) in Theorem 3.1 is important. It also confirms that one cannot omit the condition $\sup_{(x,a) \in \mathbf{X} \times \mathbf{A}} |c(x, a)| < \infty$ in Proposition 9.13 of [1].

The MDP we plan to consider is discounted, with discrete state and actions spaces. In this case, the model is absorbing at Δ under each initial distribution and, traditionally, the notations are slightly modified. First of all, with some abuse of notation, we write $p(y|x, a)$ for stochastic kernels, rather than $p(\{y\}|x, a)$. Secondly, the cemetery Δ is omitted and the transition probability on $\mathbf{X} \not\rightarrow \Delta$ is normalized:

$$p(y|x, a) = \beta q(y|x, a) \text{ for } y \in \mathbf{X};$$

The initial distribution P_0 is on $\mathcal{B}(\mathbf{X})$.

Now equations (3.1) and (3.2) take the form

$$v_x^{\hat{\pi}} = \inf_{a \in \mathbf{A}} \left\{ c(x, a) + \beta \int_{\mathbf{X}} v_y^{\hat{\pi}} q(dy|x, a) \right\}, \quad x \in \mathbf{X}. \quad (4.1)$$

and

$$\liminf_{T \rightarrow \infty} \beta^T E_{P_0}^\pi [v_{X_T}^{\hat{\pi}}] \leq 0: \quad (4.2)$$

in the original notations, $(1 - \beta^T)$ is the probability that $X_T = \Delta$ leading to $v_{X_T}^{\hat{\pi}} = 0$.

Let $\mathbf{X} = \{0, 1, 2, \dots\}$, $\mathbf{A} = \{1, 2\}$, $q(0|0, a) \equiv 1$, $c(0, a) \equiv 0$. For $x > 0$ we put $q(x+1|x, 1) = q(0|x, 1) \equiv 1/2$, $q(x+1|x, 2) \equiv 1$, other transition probabilities $q(\cdot)$ being zero; $c(x, 1) = 2^x$, $c(x, 2) \equiv 1$. The discount factor is $\beta = 1/2$. (See Fig. 1, where the probabilities q are shown.)

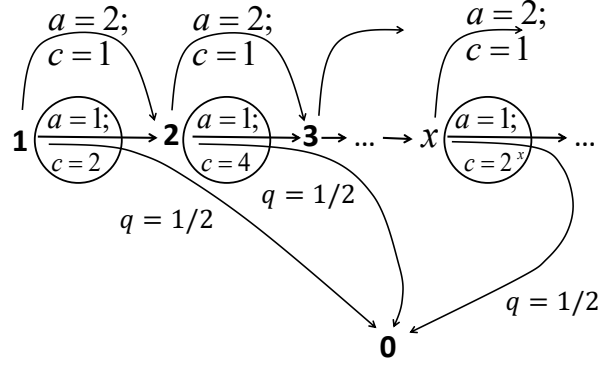


FIGURE 1. The selector $\hat{\varphi}(x) \equiv 1$ is not optimal.

The optimality equation

$$v(x) = \inf_{a \in \mathbf{A}} \left\{ c(x, a) + \beta \int_{\mathbf{X}} v(y) q(dy|x, a) \right\}, \quad x \in \mathbf{X} \quad (4.3)$$

takes the form

$$v(0) = \frac{1}{2}v(0);$$

$$\text{for } x > 0, \quad v(x) = \min\left\{2^x + \frac{1}{4}v(x+1) + \frac{1}{4}v(0); 1 + \frac{1}{2}v(x+1)\right\},$$

and has the minimal non-negative solution

$$v(x) = v_x^* \equiv 2 \text{ for } x > 0; \quad v(0) = 0,$$

coincident with the Bellman function by [10, Thm. 7.3.2]. The stationary selector $\varphi^*(x) \equiv 2$ is conserving and equalizing and hence optimal for each initial state $x \in \mathbf{X}$: see [1, Prop. 9.12].

Now consider the stationary selector $\hat{\varphi}(x) \equiv 1$. The performance functional $v_x^{\hat{\varphi}}$ is given by $v_x^{\hat{\varphi}} = 2^{x+1}$ (for $x > 0$) and satisfies equation (4.1):

$$1 + \frac{1}{2}v_{x+1}^{\hat{\varphi}} = 1 + 2^{x+1} > 2^{x+1} = v_x^{\hat{\varphi}} = 2^x + \frac{1}{4}v_{x+1}^{\hat{\varphi}}.$$

Equation (4.2) is violated for selector φ^* :

$$\beta^T \mathbb{E}_x^{\varphi^*} \left[v_{X_T}^{\hat{\varphi}} \right] = \beta^T 2^{x+T+1} = 2^{x+1},$$

and the selector $\hat{\varphi}$ is certainly non-optimal for each initial distribution P_0 , except $P_0(dx) = \delta_0(dx)$. If we put $P_0(dx) = \delta_{x_0}(dx)$ for an arbitrarily fixed $x_0 > 0$, we see that all the conditions in Theorem 3.1 are satisfied for $\hat{\varphi}$ apart from equation (3.2). In particular, for each strategy π , $\mathbb{E}_x^\pi \left[v_{X_t}^{\hat{\varphi}} \right] \leq \mathbb{E}_x^{\varphi^*} \left[v_{X_t}^{\hat{\varphi}} \right] = 2^{x+t+1}$ is finite: the probability to reach state $X_T \neq 0$ is the highest for φ^* . By the way,

$$\beta^T \mathbb{E}_x^{\hat{\varphi}} \left[v_{X_T}^{\hat{\varphi}} \right] = 2^{x+1-T} \rightarrow 0 \text{ as } T \rightarrow \infty.$$

We see that the both functions $v_x^* = v_x^{\varphi^*}$ and $v_x^{\hat{\varphi}}$ solve the optimality equation (4.3), but only v_x^* is the minimal non-negative solution. Note that the optimality equation (4.3) has many other solutions, e.g. $v(x) = 2 + k \cdot 2^x$ with $k \in [0, 1/2]$.

5. CONCLUSION

Theorem 3.1 makes it possible to prove optimality of a strategy without knowing the Bellman function. Similar statements in the existing literature are known under additional restrictive conditions. At the same time, Theorem 3.1 can also be used for calculating the Bellman function. Indeed, if that theorem is applicable for each degenerate initial distribution $P_0(dy) = \delta_x(dy)$, $x \in \mathbf{X}$ (the Dirac measure), then

$$\mathbb{E}_x^{\hat{\pi}} \left[\sum_{t=1}^{\infty} c(X_{t-1}, A_t) \right] = v_x^{\hat{\pi}} = v_x^*$$

equals the Bellman function. In this connection, let us remind that the Bellman equation

$$v_x^* = \inf_{a \in \mathbf{A}} \left\{ c(x, a) + \int_{\mathbf{X}} v_y^* p(dy|x, a) \right\}, x \in \mathbf{X}$$

in an infinite-horizon model can have many solutions even in finite models [3, Ch. 4, Section 7], [9, Subsection 2.2.3], [10, Ex. 7.2.3 and 7.3.1], see also the above counter-example. Sometimes it is hard to choose the Bellman function among those solutions.

REFERENCES

- [1] D. Bertsekas, S. Shreve, Stochastic Optimal Control, Academic Press, New York, 1978.
- [2] D. Bertsekas, Dynamic Programming and Optimal Control, V.II, Athena Scientific, Belmont, MA, USA, 2001.
- [3] E.B. Dynkin, A.A. Yushkevich, Controlled Markov Processes, Springer, New York, 1979.
- [4] E.A. Feinberg, U. Rothblum, Splitting randomized stationary policies in total-reward Markov decision processes, Math. Oper. Res. 37 (2012) 129-153.
- [5] P.G. Canbolat, U.G. Rothblum, (Approximate) iterated successive approximations algorithm for sequential decision processes, Ann. Oper. Res. 208 (2013) 309-320.
- [6] O. Hernández-Lerma, J. Lasserre, Discrete-Time Markov Control Processes, Springer Verlag, New York, 1996.
- [7] Hernández-Lerma, J. Lasserre, Further Topics in Discrete-Time Markov Control Processes, Springer Verlag, New York, 1999.
- [8] A. Piunovskiy, Optimal Control of Random Sequences in Problems with Constraints, Kluwer, Dordrecht, 1997.
- [9] A.B. Piunovskiy, Examples in Markov Decision Processes, Imperial College Press, London, 2013.
- [10] M.L. Puterman, Markov Decision Processes, Wiley, New York, 1994.