# RE-EXAMINATION OF MARKOV POLICIES FOR ADDITIVE DECISION PROCESS

Fujita, Toshiharu
Graduate School of Mathematics, Kyushu University

# RE-EXAMINATION OF MARKOV POLICIES FOR ADDITIVE DECISION PROCESS

By

Toshiharu FUJITA*

### Abstract

The purpose of this paper is to ensure that Markov policy is enough for the additive decision process. We show that there exists an optimal policy which is Markov for both stochastic additive decision process and deterministic one. We also verify through a multi-stage stochastic decision tree method that, among the class of general policies, there exist an optimal Markov policy. This fact is of course obtained by solving the regular recursive equation.

## 1. Introduction

In this paper we consider a class of sequential optimization problems whose objective function is additive. We focus our attention on the tacitly known but never clearly proved fact that Markov policy is enough in the so-called Markov decision processes (Bellman and Zadeh(1970, p. 152, *ll.* 19-22), Bertsekas and Shreve(1978, p. 6, *ll.* 20-23), Howard(1960), Stokey and Lucas Jr.(1989), Puterman(1994), Iwamoto, Tsurusaki and Fujita, and others). It has long been thought that there is no room for argument on this point. However, we must accept a recent result that an optimal policy is not Markov for a certain problem especially for minimum objective function problem (Iwamoto and Fujita(1995), Iwamoto, Tsurusaki and Fujita). These two apparently inconsistent facts enable us to reconsider whether Markov policy is enough or not.

In section 2 we discuss the stochastic maximization of additive function. We derive a recursive equation both for the class of general policies and for the class of Markov policies. Verifying that the optimal value functions in both classes are identical, we show that Markov policy is enough. In section 3 we give rigorous proofs of theorems in section 2. In section 4 we show the corresponding results for the deterministic optimization, which is a degenerate case of the stochastic one. In section 5, illustrating multi-stage stochastic decision trees, we give a two-stage stochastic decision process, a pair of optimal value functions, and an optimal Markov policy. The pair is also obtained by solving the corresponding recursive equation.

* Graduate School of Mathematics, Kyushu University 33, Fukuoka 812-81, Japan.
  Research Fellow of the Japan Society for the Promotion of the Science.
  e-mail: fujita@math.kyushu-u.ac.jp

Throughout the paper the following definition and notation are used :

$N \geq 2$ is an integer ; the *total number of stages*

$X = \{s_1,\, s_2, \ldots, s_p\}$ is a finite *state space*

$U = \{a_1,\, a_2, \ldots, a_k\}$ is a finite *action space*

$r_n : X \times U \to R^1$ is an *n-th reward function*   $(1 \leq n \leq N)$

$r_G : X \to R^1$ is a *terminal reward function*

$f : X \times U \to X$ is a *deterministic transition law*

   ; $f(x, u)$ represents the successor state of $x$ for action $u$

$p$ is a *Markov transition law*

   : $p(y|x, u) \geq 0$  $\forall (x, u, y) \in X \times U \times X$,   $\displaystyle\sum_{y \in X} p(y|x, u) = 1$  $\forall (x, u) \in X \times U$

$y \sim p(\,\cdot\,|x, u)$ denotes that next state $y$ conditioned on state $x$ and action $u$ appears with probability $p(y|x, u)$.

## 2.   Stochastic Maximization

We consider the stochastic maximization problem with additive function as follows :

$$\text{Maximize} \quad E[\, r_1(x_1, u_1) + r_2(x_2, u_2) + \cdots + r_N(x_N, u_N) + r_G(x_{N+1})\,]$$

$$\text{subject to} \quad \text{(i)} \quad x_{n+1} \sim p(\,\cdot\,|x_n, u_n) \tag{1}$$

$$\text{(ii)} \quad u_n \in U \qquad\qquad n = 1, 2, \ldots, N$$

where the sequence of states $\{x_2,\, x_3,\, \ldots,\, x_{N+1}\}$ together with a sequence of intermediate actions $\{u_1,\, u_2,\, \ldots,\, u_N\}$ is stochastically generated through an initial state $x_1$, the Markov transition law $x_{n+1} \sim p(\,\cdot\,|x_n, u_n)$ and a (general or Markov) policy.

### 2.1. General Policies

First we consider the problem (1) with the set of all general policies. In general, a *policy* $\sigma = \{\sigma_1,\, \sigma_2,\, \ldots,\, \sigma_N\}$ is a sequence of *decision functions*

$$\sigma_1 : X \to U, \quad \sigma_2 : X \times X \to U, \quad \ldots \quad, \quad \sigma_N : \underbrace{X \times \cdots \times X}_{N \text{ times}} \to U. \tag{2}$$

In what follows, in order to distinguish the policy, we call this policy *general* policy. Note that the size in our data specified in (2) yields $k^{p^n}$ *n-th* decision functions $\sigma_n$ ($n = 1, 2, \ldots, N$) and $k^{p + p^2 + \cdots + p^N}$ general policies $\sigma$.

An application of general policy $\sigma$ at an initial state $x_1$ generates stochastically the alternate sequence of states and actions $\{u_1, x_2, u_2, x_3, \ldots, u_N, x_{N+1}\}$ as follows :

$$
\begin{aligned}
\sigma_1(x_1) &= u_1 & &\to & p(\,\cdot\,|x_1, u_1) &\sim x_2 & &\to \\
\sigma_2(x_1, x_2) &= u_2 & &\to & p(\,\cdot\,|x_2, u_2) &\sim x_3 & &\to \\
&\ \vdots & & & &\ \vdots \\
\sigma_N(x_1, x_2, \ldots, x_N) &= u_N & &\to & p(\,\cdot\,|x_N, u_N) &\sim x_{N+1}.
\end{aligned}
$$

We call this problem *general problem*. With any general policy $\sigma = \{\sigma_n, \ldots, \sigma_N\}$ over the $(N - n + 1)$-stage process starting on $n$-th stage and terminating at the last stage, we associate the expected value :

$$I^n(x_n; \sigma) = \sum_{(x_{n+1}, \ldots, x_N) \in X \times \cdots \times X} \sum \cdots \sum \{[r_n(x_n, u_n) + \cdots + r_N(x_N, u_N) + r_G(x_{N+1})]$$
$$\times p(x_{n+1}|x_n, u_n) \cdots p(x_{N+1}|x_N, u_N)\} \quad (3)$$

where $\{u_n, x_{n+1}, u_{n+1}, x_{n+2}, \ldots, u_N, x_{N+1}\}$ is stochastically generated through the general policy $\sigma$ and the starting state $x_n$ as follows :

$$
\begin{array}{llll}
\sigma_n(x_n) = u_n & \rightarrow & p(\cdot|x_n, u_n) \sim x_{n+1} & \rightarrow \\
\sigma_{n+1}(x_n, x_{n+1}) = u_{n+1} & \rightarrow & p(\cdot|x_{n+1}, u_{n+1}) \sim x_{n+2} & \rightarrow \\
\quad \vdots & & \quad \vdots & \\
\sigma_N(x_n, x_{n+1}, \ldots, x_N) = u_N & \rightarrow & p(\cdot|x_N, u_N) \sim x_{N+1}.
\end{array}
\quad (4)
$$

We define the family of the corresponding *general subproblems* as follows :

$$
\begin{aligned}
V^{N+1}(x_{N+1}) &= r_G(x_{N+1}) & x_{N+1} \in X \\
V^n(x_n) &= \operatorname*{Max}_{\sigma} I^n(x_n; \sigma) & x_n \in X, \; 1 \le n \le N.
\end{aligned}
\quad (5)
$$

Note that the general problem (1) is identical to (5) with $n = 0$. Furthermore we should remark that the maximization for the subproblems stated above is taken for all general policies, namely, in problem (5)

$$\sigma_n : X \rightarrow U, \quad \sigma_{n+1} : X \times X \rightarrow U, \quad \ldots, \quad \sigma_N : \underbrace{X \times \cdots \times X}_{N-n+1 \text{ times}} \rightarrow U.$$

Then we have the recursive formula for the general subproblems :

THEOREM 2.1.

$$
\begin{aligned}
V^{N+1}(x) &= r_G(x) & x \in X \\
V^n(x) &= \operatorname*{Max}_{u \in U}[r_n(x, u) + \sum_{y \in X} V^{n+1}(y)p(y|x, u)] & x \in X, \; 1 \le n \le N.
\end{aligned}
\quad (6)
$$

## 2.2. Markov Policies

In this subsection we restrict the problem (1) to the set of all Markov policies. We call this problem *Markov problem*. Here a policy

$$\pi = \{\pi_1, \pi_2, \ldots, \pi_N\}$$

is called *Markov* if

$$\pi_1 : X \rightarrow U, \quad \pi_2 : X \rightarrow U, \quad \ldots, \quad \pi_N : X \rightarrow U. \quad (7)$$

Thus, any Markov policy $\pi$ with an initial state $x_1$ yields the Markov chain on $X$ and the sequence of the resulting actions as follows :

$$
\begin{aligned}
\pi_1(x_1) = u_1 &\quad \rightarrow \quad p(\cdot|x_1, u_1) \sim x_2 &\quad \rightarrow \\
\pi_2(x_2) = u_2 &\quad \rightarrow \quad p(\cdot|x_2, u_2) \sim x_3 &\quad \rightarrow \\
\vdots &\qquad\qquad \vdots \\
\pi_N(x_N) = u_N &\quad \rightarrow \quad p(\cdot|x_N, u_N) \sim x_{N+1}.
\end{aligned}
$$

We remark that the size in (7) yields $k^p$ $n$-th decision functions $\pi_n$ $(n = 1, 2, \ldots, N)$ and $k^{Np}$ Markov policies $\pi$.

Note that any Markov policy $\pi = \{\pi_n, \ldots, \pi_N\}$ over the $(N-n+1)$-stage process is associated with its expected value $I^n(x_n; \pi)$ defined by (3), where the alternate sequence $\{u_n, x_{n+1}, u_{n+1}, x_{n+2}, \ldots, u_N, x_{N+1}\}$ is similarly generated as in (4). Here we remark that

$$
u_n = \pi_n(x_n), \quad u_{n+1} = \pi_{n+1}(x_{n+1}), \quad \cdots, u_N = \pi_N(x_N).
$$

We define the corresponding *Markov subproblems* as follows :

$$
\begin{aligned}
v^{N+1}(x_{N+1}) &= r_G(x_{N+1}) &\quad x_{N+1} \in X \\
v^n(x_n) &= \operatorname*{Max}_{\pi} I^n(x_n; \pi) &\quad x_n \in X, \ 1 \le n \le N.
\end{aligned} \tag{8}
$$

Then (8) with $n = 1$ reduces to the Markov problem (1). We have the recursive formula for the Markov subproblems :

**THEOREM 2.2.**

$$
\begin{aligned}
v^{N+1}(x) &= r_G(x) &\quad x \in X \\
v^n(x) &= \operatorname*{Max}_{u \in U}[r_n(x, u) + \sum_{y \in X} v^{n+1}(y) p(y|x, u)] &\quad x \in X, \ 1 \le n \le N.
\end{aligned} \tag{9}
$$

**THEOREM 2.3.** (i) *A Markov policy yields the optimal value function $V^1(\cdot)$ for the general problem. That is, there exists an optimal Markov policy $\pi^*$ for the general problem (1) :*

$$
I^1(x_1; \pi^*) = V^1(x_1) \quad \text{for all} \ x_1 \in X.
$$

*In fact, letting $\pi_n^*(x)$ be a maximizer of (9) (or (6)) for each $x \in X$, $1 \le n \le N$, we have the optimal Markov policy $\pi^* = \{\pi_1^*, \ldots, \pi_N^*\}$.*

(ii) *The optimal value functions for the Markov subproblems (8) are equal to the optimal value functions for the general subproblems (5) :*

$$
v^n(x) = V^n(x) \quad x \in X, \ 1 \le n \le N + 1.
$$

## 3.   Proof of Theorems

In this subsection we prove Theorems 2.1 - 2.3. We remark that Theorem 2.3 (i) implies Theorem 2.3 (ii) and that a combination of Theorem 2.1 and Theorem 2.3 (ii)

yields Theorem 2.2. Thus it suffices to prove Theorem 2.1 and Theorem 2.3 (i). Since there is no essential difficulty in extending the proof to the general $N$-stage process, we prove both theorems for the two-stage process, namely, for the case $N = 2$.

We note that

$$
\begin{aligned}
V^3(x_3) &= r_G(x_3) \\
V^2(x_2) &= \underset{\sigma_2}{\text{Max}} \sum_{x_3 \in X} [\, r_2(x_2, u_2) + r_G(x_3)\,] p(x_3|x_2, u_2) \quad &(10) \\
V^1(x_1) &= \underset{\sigma_1, \sigma_2}{\text{Max}} \sum_{(x_2, x_3) \in X \times X} \sum \{[\, r_1(x_1, u_1) + r_2(x_2, u_2) + r_G(x_3)\,] \quad &(11)
\end{aligned}
$$

$$
\times p(x_2|x_1, u_1) p(x_3|x_2, u_2)\}
$$

where $u_2 = \sigma_2(x_2)$ in (10) and $u_1 = \sigma_1(x_1)$, $u_2 = \sigma_2(x_1, x_2)$ in (11), respectively.

Thus the equality

$$
V^2(x_2) = \underset{u_2 \in U}{\text{Max}}[r_2(x_2, u_2) + \sum_{x_3 \in X} V^3(x_3) p(x_3|x_2, u_2)] \quad x_2 \in X
$$

is trivial. We prove

$$
V^1(x_1) = \underset{u_1 \in U}{\text{Max}}[r_1(x_1, u_1) + \sum_{x_2 \in X} V^2(x_2) p(x_2|x_1, u_1)] \quad x_1 \in X. \quad (12)
$$

Let us choose an optimal (necessarily Markov) policy $\sigma_2^*$ for the one-stage process :

$$
V^2(x_2) = I^2(x_2; \sigma_2^*) \quad \forall x_2 \in X. \quad (13)
$$

From the definition (5), we can for each $x_1 \in X$ choose an optimal (not necessarily Markov) policy $\tilde{\sigma} = \{\tilde{\sigma}_1, \tilde{\sigma}_2\}$ for the two-stage process :

$$
V^1(x_1) = I^1(x_1; \tilde{\sigma}) \quad x_1 \in X.
$$

Thus we see that

$$
\begin{aligned}
&V^1(x_1) \\
=\ & I^1(x_1; \tilde{\sigma}_1, \tilde{\sigma}_2) \\
=\ & \sum_{(x_2, x_3) \in X \times X} \sum \{[r_1(x_1, u_1) + r_2(x_2, u_2) + r_G(x_3)] p(x_2|x_1, u_1) p(x_3|x_2, u_2)\} \quad (14)
\end{aligned}
$$

where

$$
u_1 = \tilde{\sigma}_1(x_1), \ u_2 = \tilde{\sigma}_2(x_1, x_2).
$$

Since

$$
\begin{aligned}
&\sum_{(x_2, x_3) \in X \times X} \sum \{[r_1(x_1, u_1) + r_2(x_2, u_2) + r_G(x_3)] p(x_2|x_1, u_1) p(x_3|x_2, u_2)\} \\
=\ & \sum_{x_2 \in X} \{r_1(x_1, u_1) + \sum_{x_3 \in X} [r_2(x_2, u_2) + r_G(x_3)] p(x_3|x_2, u_2)\} p(x_2|x_1, u_1)
\end{aligned}
$$

and

$$\sum_{x_3 \in X} [r_2(x_2, u_2) + r_G(x_3)] p(x_3 | x_2, u_2) \leq I^2(x_2; \sigma_2^*) = V^2(x_2) \quad \forall x_2 \in X,$$

we have from (14)

$$\begin{aligned}
V^1(x_1) &\leq \sum_{x_2 \in X} [r_1(x_1, u_1) + V^2(x_2)] p(x_2 | x_1, u_1) \\
&= r_1(x_1, u_1) + \sum_{x_2 \in X} V^2(x_2) p(x_2 | x_1, u_1).
\end{aligned}$$

Thus taking maximum over $u \in U$, we get

$$V^1(x_1) \leq \underset{u_1 \in U}{\text{Max}} [r_1(x_1, u_1) + \sum_{x_2 \in X} V^2(x_2) p(x_2 | x_1, u_1)] \quad \forall x_1 \in X. \tag{15}$$

On the other hand, let for any $x_1 \in X$, $u^* = u^*(x_1) \in U$ be a maximizer of the right hand side of (15). This defines a Markov decision function

$$\pi_1^* : X \to U \quad \pi_1^*(x_1) = u^*(x_1).$$

Then we have

$$\begin{aligned}
&\underset{u_1 \in U}{\text{Max}} [r_1(x_1, u_1) + \sum_{x_2 \in X} V^2(x_2) p(x_2 | x_1, u_1)] \\
&= r_1(x_1, u_1) + \sum_{x_2 \in X} V^2(x_2) p(x_2 | x_1, u_1) \quad (u_1 = \pi_1^*(x_1)). \tag{16}
\end{aligned}$$

From (13), we get

$$V^2(x_2) = \sum_{x_3 \in X} [r_2(x_2, u_2) + r_G(x_3)] p(x_3 | x_2, u_2) \quad (u_2 = \sigma_2^*(x_2)). \tag{17}$$

Thus we have from (17)

$$\begin{aligned}
&r_1(x_1, u_1) + \sum_{x_2 \in X} V^2(x_2) p(x_2 | x_1, u_1) \quad (u_1 = \pi_1^*(x_1)) \\
&= r_1(x_1, u_1) + \sum_{x_2 \in X} [\sum_{x_3 \in X} [r_2(x_2, u_2) + r_G(x_3)] p(x_3 | x_2, u_2)] p(x_2 | x_1, u_1) \\
&= \sum_{(x_2, x_3) \in X \times X} \{[r_1(x_1, u_1) + r_2(x_2, u_2) + r_G(x_3)] p(x_2 | x_1, u_1) p(x_3 | x_2, u_2)\}. \tag{18}
\end{aligned}$$

Combining (16) and (18), we obtain

$$\begin{aligned}
&\underset{u_1 \in U}{\text{Max}} [r_1(x_1, u_1) + \sum_{x_2 \in X} V^2(x_2) p(x_2 | x_1, u_1)] \\
&= \sum_{(x_2, x_3) \in X \times X} \{[r_1(x_1, u_1) + r_2(x_2, u_2) + r_G(x_3)] p(x_2 | x_1, u_1) p(x_3 | x_2, u_2)\}
\end{aligned}$$

$$(u_1 = \pi_1^*(x_1), \ u_2 = \sigma_2^*(x_2))$$

$$\le \ \underset{\sigma_1, \sigma_2}{\text{Max}} \sum \sum_{(x_2, x_3) \in X \times X} \{[r_1(x_1, u_1) + r_2(x_2, u_2) + r_G(x_3)]$$

$$\times p(x_2 | x_1, u_1) p(x_3 | x_2, u_2)\}$$

$$= \ V^1(x_1). \tag{19}$$

Both equations (15) and (19) imply the desired equality (12). This completes the proof of Theorem 2.1.

Furthermore, the equalities in (19) imply that the optimal value function $V^1(\cdot)$ is yielded by the Markov policy $\bar{\pi} = \{\pi_1^*, \sigma_2^*\}$ :

$$V^1(x_1) = I^1(x_1; \bar{\pi}) \quad x_1 \in X.$$

This completes the proof of Theorem 2.3 (i).

## 4. Deterministic Maximization

In this section we consider the deterministic maximization problem with additive function as follows :

$$\text{Maximize} \quad r_1(x_1, u_1) + r_2(x_2, u_2) + \cdots + r_N(x_N, u_N) + r_G(x_{N+1})$$

$$\text{subject to} \quad \text{(i)} \quad f(x_n, u_n) = x_{n+1} \tag{20}$$

$$\text{(ii)} \quad u_n \in U \qquad n = 1, 2, \ldots, N.$$

Note that this problem is the special case of the stochastic maximization problem (1). Through this section the problem (20) with the set of all general (resp. Markov) policies

$$\sigma = \{\sigma_1, \sigma_2, \ldots, \sigma_N\} \quad (\text{resp. } \pi = \{\pi_1, \pi_2, \ldots, \pi_N\})$$

is called the general (resp. Markov) problem.

First we consider the general problem. We associate any policy $\sigma = \{\sigma_n, \ldots, \sigma_N\}$ for the $(N - n + 1)$-stage process starting on $n$-th stage and terminating at the last stage with its value :

$$I^n(x_n; \sigma) = r_n(x_n, u_n) + \cdots + r_N(x_N, u_N) + r_G(x_{N+1}) \tag{21}$$

where $\{u_n, x_{n+1}, u_{n+1}, x_{n+2}, \ldots, u_N, x_{N+1}\}$ is uniquely determined through $\sigma$ and $x_n$ as follows :

$$\begin{aligned} \sigma_n(x_n) = u_n & \quad \rightarrow \quad f(x_n, u_n) = x_{n+1} & \quad \rightarrow \\ \sigma_{n+1}(x_n, x_{n+1}) = u_{n+1} & \quad \rightarrow \quad f(x_{n+1}, u_{n+1}) = x_{n+2} & \quad \rightarrow \\ \vdots & \qquad\qquad \vdots & \\ \sigma_N(x_1, x_1, \ldots, x_N) = u_N & \quad \rightarrow \quad f(x_N, u_N) = x_{N+1}. & \end{aligned} \tag{22}$$

We consider the following family of subproblems :

$$V^{N+1}(x_{N+1}) = r_G(x_{N+1}) \qquad x_{N+1} \in X$$

$$V^n(x_n) = \underset{\sigma}{\text{Max}}\, I^n(x_n; \sigma) \qquad x_n \in X, \ 1 \le n \le N. \tag{23}$$

Note that the general problem is identical to (23) with $n = 0$. Then we have the backward recursive formula :

THEOREM 4.1.

$$V^{N+1}(x) = r_G(x) \qquad x \in X$$
$$V^n(x) = \underset{u \in U}{\text{Max}}[r_n(x, u) + V^{n+1}(f(x, u))] \qquad x \in X, \quad 1 \le n \le N. \qquad (24)$$

Next we restrict the problem (20) to the set of all Markov policies. Note that any Markov policy $\pi = \{\pi_n, \dots, \pi_N\}$ over the $(N-n+1)$-stage process is associated with its value $I^n(x_n; \pi)$ defined by (21), where the alternate sequence $\{u_n, x_{n+1}, u_{n+1}, x_{n+2}, \dots , u_N, x_{N+1}\}$ is similarly determined through the Markov policy $\pi$ and the starting state $x_n$ as in (22).

We define the corresponding *Markov subproblems* as follows :

$$v^{N+1}(x_{N+1}) = r_G(x_{N+1}) \qquad x_{N+1} \in X$$
$$v^n(x_n) = \underset{\pi}{\text{Max}}\, I^n(x_n; \pi) \qquad x_n \in X, \quad 1 \le n \le N. \qquad (25)$$

Then (25) with $n = 1$ reduces to the Markov problem. We have the recursive formula for the Markov subproblems :

THEOREM 4.2.

$$v^{N+1}(x) = r_G(x) \qquad x \in X$$
$$v^n(x) = \underset{u \in U}{\text{Max}}[r_n(x, u) + v^{n+1}(f(x, u))] \qquad x \in X, \quad 1 \le n \le N. \qquad (26)$$

THEOREM 4.3. (i) *A Markov policy yields the optimal value function* $V^1(\cdot)$ *for the general problem. That is, there exists an optimal Markov policy* $\pi^*$ *for the general problem (20) :*

$$I^1(x_1; \pi^*) = V^1(x_1) \quad \text{for all} \ \ x_1 \in X.$$

*In fact, letting* $\pi_n^*(x)$ *be a maximizer of (26) (or (24)) for each* $x \in X$, $1 \le n \le N$, *we have the optimal Markov policy* $\pi^* = \{\pi_1^*, \dots, \pi_N^*\}$.

(ii) *The optimal value functions for the Markov subproblems (25) are equal to the optimal value functions for the general subproblems (23) :*

$$v^n(x) = V^n(x) \qquad x \in X, \quad 1 \le n \le N + 1.$$

Since Theorems 4.1 - 4.3 are special cases of Theorems 2.1 - 2.3 respectively, each theorem in this section is clear.

## 5.   Example

We illustrate the following two-stage, three-state and two-action stochastic decision process :

$$\text{Maximize} \quad E[r_1(u_1) + r_2(u_2) + r_G(x_3)]$$
$$\text{subject to} \quad \text{(i)} \ \ x_{n+1} \sim p(\cdot | x_n, u_n) \quad n = 1, 2 \qquad (27)$$
$$\text{(ii)} \ \ u_1 \in U, \ u_2 \in U$$

where the data is

$$r_G(s_1) = 0.3 \qquad r_G(s_2) = 1.0 \qquad r_G(s_3) = 0.8 \tag{28}$$

$$r_2(a_1) = 1.0 \qquad r_2(a_2) = 0.6$$

$$r_1(a_1) = 0.7 \qquad r_1(a_2) = 1.0$$

$u_t = a_1$

| $x_t \backslash x_{t+1}$ | $s_1$ | $s_2$ | $s_3$ |
|---|---|---|---|
| $s_1$ | 0.8 | 0.1 | 0.1 |
| $s_2$ | 0.0 | 0.1 | 0.9 |
| $s_3$ | 0.8 | 0.1 | 0.1 |

$u_t = a_2$

| $x_t, \backslash x_{t+1}$ | $s_1$ | $s_2$ | $s_3$ |
|---|---|---|---|
| $s_1$ | 0.1 | 0.9 | 0.0 |
| $s_2$ | 0.8 | 0.1 | 0.1 |
| $s_3$ | 0.1 | 0.0 | 0.9 |

In order to solve the problem, we directly generate one- and two- stage stochastic decision trees and enumerate all the possible histories together with the related expected values. We call this brute force enumeration method a *multi-stage stochastic decision tree method*. For any given policy, this tree method traces all the resulting histories. Then it yields the value of the policy. Further, from among all general policies, it selects an optimal policy together with the sequence of optimal value functions. The multi-stage stochastic decision tree method is also applies to non-additive problems (Iwamoto and Fujita(1995), Iwamoto, Tsurusaki and Fujita).

We remark that the size yields $2^3 = 8$ first decision functions $\sigma_1 = \begin{pmatrix} \sigma_1(s_1) \\ \sigma_1(s_2) \\ \sigma_1(s_3) \end{pmatrix}$

and $2^{3 \times 3} = 512$ second decision functions

$$\sigma_2 = \begin{pmatrix} \sigma_2(s_1, s_1) & \sigma_2(s_2, s_1) & \sigma_2(s_3, s_1) \\ \sigma_2(s_1, s_2) & \sigma_2(s_2, s_2) & \sigma_2(s_3, s_2) \\ \sigma_2(s_1, s_3) & \sigma_2(s_2, s_3) & \sigma_2(s_3, s_3) \end{pmatrix}.$$

As a total, there are $8 \times 512 = 4096$ general policies $\sigma = \{\sigma_1, \sigma_2\}$ for the problem (27). First, we have from definition (28)

$$V^3(s_1) = 0.3, \quad V^3(s_2) = 1.0, \quad V^3(s_3) = 0.8. \tag{29}$$

Second, the decision tree in Figure 1 shows

$$V^2(s_1) = 1.53, \quad V^2(s_2) = 1.82, \quad V^2(s_3) = 1.42. \tag{30}$$

Third, the enumeration in Figures 2, 3 and 4 calculates the maximum expected values :

$$V^1(s_1) = 2.791, \quad V^1(s_2) = 2.548, \quad V^1(s_3) = 2.431. \tag{31}$$

The calculation yields, at the same time, the optimal policy $\sigma^* = \{\sigma_1^*(x_1), \sigma_2^*(x_1, x_2)\}$ :

$$\sigma_1^*(s_1) = a_2, \quad \sigma_1^*(s_2) = a_2, \quad \sigma_1^*(s_3) = a_2$$

$$\sigma_2^*(s_1, s_1) = a_2, \qquad \sigma_2^*(s_2, s_1) = a_2, \qquad \sigma_2^*(s_3, s_1) = a_2$$
$$\sigma_2^*(s_1, s_2) = a_1, \qquad \sigma_2^*(s_2, s_2) = a_1 \text{ or } a_2, \quad \sigma_2^*(s_3, s_2) = a_1$$
$$\sigma_2^*(s_1, s_3) = a_1 \text{ or } a_2, \quad \sigma_2^*(s_2, s_3) = a_1, \qquad \sigma_2^*(s_3, s_3) = a_1.$$

Thus the general policy $\sigma^*$ reduces a Markov policy $\pi^* = \{\pi_1^*(x_1), \ \pi_2^*(x_2)\}$ :

$$\pi_1^*(s_1) = a_2, \quad \pi_1^*(s_2) = a_2, \quad \pi_1^*(s_3) = a_2$$

$$\pi_2^*(s_1) = a_2, \quad \pi_2^*(s_2) = a_1, \quad \pi_2^*(s_3) = a_1.$$

Thus, the Markov policy $\pi^*$ is optimal. Finally we remark that the pair of optimal value functions (29),(30),(31) and the optimal Markov policy $\pi^*$ is also obtained by solving either the corresponding recursive equation (6) or (9). Solving the latter is the so-called dynamic programming method.

Figure 1 : all one-stage behaviors from $s_1, s_2$ and $s_3$, and selection of maximum branch

$$\left( V^2(x_2) = \underset{u_2}{\text{Max}} \sum_{x_3 \in X} [\, r_2(u_2) + r_G(x_3)\,] p(x_3|x_2, u_2) \quad x_2 = s_1, \ s_2, \ s_3 \right)$$

| history | | ter. | path | sum | mult. | exp. |
|---|---|---|---|---|---|---|
| | $s_1$ | 0.3 | 0.8 | 1.3 | 1.04 | |
| | $s_2$ | 1.0 | 0.1 | 2.0 | 0.2 | 1.42 |
| | $s_3$ | 0.8 | 0.1 | 1.8 | 0.18 | |
| | $s_1$ | 0.3 | 0.1 | 0.9 | 0.09 | |
| | $s_2$ | 1.0 | 0.9 | 1.6 | 1.44 | **1.53** |
| | $s_3$ | 0.8 | 0.0 | 1.4 | 0 | |
| | $s_1$ | 0.3 | 0.0 | 1.3 | 0 | |
| | $s_2$ | 1.0 | 0.1 | 2.0 | 0.2 | **1.82** |
| | $s_3$ | 0.8 | 0.9 | 1.8 | 1.62 | |
| | $s_1$ | 0.3 | 0.8 | 0.9 | 0.72 | |
| | $s_2$ | 1.0 | 0.1 | 1.6 | 0.16 | 1.02 |
| | $s_3$ | 0.8 | 0.1 | 1.4 | 0.14 | |
| | $s_1$ | 0.3 | 0.8 | 1.3 | 1.04 | |
| | $s_2$ | 1.0 | 0.1 | 2.0 | 0.2 | **1.42** |
| | $s_3$ | 0.8 | 0.1 | 1.8 | 0.18 | |
| | $s_1$ | 0.3 | 0.1 | 0.9 | 0.09 | |
| | $s_2$ | 1.0 | 0.0 | 1.6 | 0 | 1.35 |
| | $s_3$ | 0.8 | 0.9 | 1.4 | 1.26 | |

In Figure 1 we use the following list of simplified notations :

history $= x_2 \quad r_2(u_2) \, / \, u_2 \quad p(x_3 \mid x_2, u_2) \quad x_3$

ter. = terminal value $= r_G(x_3)$

path = path probability $= p(x_3 \mid x_2, u_2)$

sum = sum of the two $= r_2(u_2) + r_G(x_3)$

mult. = path $\times$ sum

exp. = expected value.

Further, the *italic* face means probability, and the *bold* face denotes a selection of maximum of up expected value or down.

Figure 2 : all two-stage behaviors from $s_1$ and selection of maximum branch

$$\left( V^1(s_1) = \max_{u_1,u_2} \sum_{(x_2,x_3)\in X\times X} \{[\,r_1(u_1)+r_2(u_2)+r_G(x_3)\,]p(x_2|s_1,u_1)p(x_3|x_2,u_2)\} \right)$$

| history | ter. | path | sum | mult. | sub. | total |
|---|---|---|---|---|---|---|
| $s_1$, $a_1$ (1.0): 0.8 $s_1$ | 0.3 | 0.64 | 2.0 | 1.28 | | |
| 0.1 $s_2$ | 1.0 | 0.08 | 2.7 | 0.216 | 1.696 | |
| 0.1 $s_3$ | 0.8 | 0.08 | 2.5 | 0.2 | | |
| $a_2$ (0.6): 0.1 $s_1$ | 0.3 | 0.08 | 1.6 | 0.128 | | |
| 0.9 $s_2$ | 1.0 | 0.72 | 2.3 | 1.656 | **1.784** | |
| 0.0 $s_3$ | 0.8 | 0.0 | 2.1 | 0 | | |
| $s_2$, $a_1$ (1.0): 0.0 $s_1$ | 0.3 | 0.0 | 2.0 | 0 | | |
| 0.1 $s_2$ | 1.0 | 0.01 | 2.7 | 0.027 | **0.252** | |
| 0.9 $s_3$ | 0.8 | 0.09 | 2.5 | 0.225 | | |
| $a_2$ (0.6): 0.8 $s_1$ | 0.3 | 0.08 | 1.6 | 0.128 | | |
| 0.1 $s_2$ | 1.0 | 0.01 | 2.3 | 0.023 | 0.172 | 2.248 |
| 0.1 $s_3$ | 0.8 | 0.01 | 2.1 | 0.021 | | |
| $s_3$, $a_1$ (1.0): 0.8 $s_1$ | 0.3 | 0.08 | 2.0 | 0.16 | | |
| 0.1 $s_2$ | 1.0 | 0.01 | 2.7 | 0.027 | **0.212** | |
| 0.1 $s_3$ | 0.8 | 0.01 | 2.5 | 0.025 | | |
| $a_2$ (0.6): 0.1 $s_1$ | 0.3 | 0.01 | 1.6 | 0.016 | | |
| 0.0 $s_2$ | 1.0 | 0.0 | 2.3 | 0 | 0.205 | |
| 0.9 $s_3$ | 0.8 | 0.09 | 2.1 | 0.189 | | |
| $s_1$, $a_1$ (1.0): 0.8 $s_1$ | 0.3 | 0.08 | 2.3 | 0.184 | | |
| 0.1 $s_2$ | 1.0 | 0.01 | 3.0 | 0.03 | 0.242 | |
| 0.1 $s_3$ | 0.8 | 0.01 | 2.8 | 0.028 | | |
| $a_2$ (0.6): 0.1 $s_1$ | 0.3 | 0.01 | 1.9 | 0.019 | | |
| 0.9 $s_2$ | 1.0 | 0.09 | 2.6 | 0.234 | **0.253** | |
| 0.0 $s_3$ | 0.8 | 0.0 | 2.4 | 0 | | |
| $s_2$, $a_1$ (1.0): 0.0 $s_1$ | 0.3 | 0.0 | 2.3 | 0 | | |
| 0.1 $s_2$ | 1.0 | 0.09 | 3.0 | 0.27 | **2.538** | |
| 0.9 $s_3$ | 0.8 | 0.81 | 2.8 | 2.268 | | **2.791** |
| $a_2$ (0.6): 0.8 $s_1$ | 0.3 | 0.72 | 1.9 | 1.368 | | |
| 0.1 $s_2$ | 1.0 | 0.09 | 2.6 | 0.234 | 1.818 | |
| 0.1 $s_3$ | 0.8 | 0.09 | 2.4 | 0.216 | | |
| $s_3$, $a_1$ (1.0): 0.8 $s_1$ | 0.3 | 0.0 | 2.3 | 0 | | |
| 0.1 $s_2$ | 1.0 | 0.0 | 3.0 | 0 | 0 | |
| 0.1 $s_3$ | 0.8 | 0.0 | 2.8 | 0 | | |
| $a_2$ (0.6): 0.1 $s_1$ | 0.3 | 0.0 | 1.9 | 0 | | |
| 0.0 $s_2$ | 1.0 | 0.0 | 2.6 | 0 | 0 | |
| 0.9 $s_3$ | 0.8 | 0.0 | 2.4 | 0 | | |

Branch structure (history tree): root $s_1$ with $a_1$ (0.7, dashed) and $a_2$ (1.0). Under $a_1$: $s_1$ (0.8), $s_2$ (0.1), $s_3$ (0.1). Under $a_2$: $s_1$ (0.1), $s_2$ (0.9), $s_3$ (0.0).

Figure 3 : all two-stage behaviors from $s_2$ and selection of maximum branch

$$\left( V^1(s_2) = \underset{u_1,u_2}{\mathrm{Max}} \sum_{(x_2,x_3)\in X\times X} \{[r_1(u_1) + r_2(u_2) + r_G(x_3)]p(x_2|s_2,u_1)p(x_3|x_2,u_2)\} \right)$$

| history | ter. | path | sum | mult. | sub. | total |
|---|---|---|---|---|---|---|
| $s_1$ (via $s_1$, $a_1$ 1.0): $s_1$ | 0.3 | 0.0 | 2.0 | 0 | | |
| $s_2$ | 1.0 | 0.0 | 2.7 | 0 | **0** | |
| $s_3$ | 0.8 | 0.0 | 2.5 | 0 | | |
| $a_2$ 0.6: $s_1$ | 0.3 | 0.0 | 1.6 | 0 | | |
| $s_2$ | 1.0 | 0.0 | 2.3 | 0 | **0** | |
| $s_3$ | 0.8 | 0.0 | 2.1 | 0 | | |
| $s_2$ (0.1) $a_1$ 1.0: $s_1$ | 0.3 | 0.0 | 2.0 | 0 | | |
| $s_2$ | 1.0 | 0.01 | 2.7 | 0.027 | **0.252** | |
| $s_3$ | 0.8 | 0.09 | 2.5 | 0.225 | | 2.16 |
| $a_2$ 0.6: $s_1$ | 0.3 | 0.08 | 1.6 | 0.128 | | |
| $s_2$ | 1.0 | 0.01 | 2.3 | 0.023 | 0.172 | |
| $s_3$ | 0.8 | 0.01 | 2.1 | 0.021 | | |
| $s_3$ (0.9) $a_1$ 1.0: $s_1$ | 0.3 | 0.72 | 2.0 | 1.44 | | |
| $s_2$ | 1.0 | 0.09 | 2.7 | 0.243 | **1.908** | |
| $s_3$ | 0.8 | 0.09 | 2.5 | 0.225 | | |
| $a_2$ 0.6: $s_1$ | 0.3 | 0.09 | 1.6 | 0.144 | | |
| $s_2$ | 1.0 | 0.0 | 2.3 | 0 | 1.845 | |
| $s_3$ | 0.8 | 0.81 | 2.1 | 1.701 | | |
| $s_1$ (0.8) $a_1$ 1.0: $s_1$ | 0.3 | 0.64 | 2.3 | 1.472 | 1.936 | |
| $s_2$ | 1.0 | 0.08 | 3.0 | 0.24 | | |
| $s_3$ | 0.8 | 0.08 | 2.8 | 0.224 | | |
| $a_2$ 0.6: $s_1$ | 0.3 | 0.08 | 1.9 | 0.152 | **2.024** | |
| $s_2$ | 1.0 | 0.72 | 2.6 | 1.872 | | |
| $s_3$ | 0.8 | 0.0 | 2.4 | 0 | | |
| $s_2$ (0.1) $a_1$ 1.0: $s_1$ | 0.3 | 0.0 | 2.3 | 0 | | |
| $s_2$ | 1.0 | 0.01 | 3.0 | 0.03 | **0.282** | |
| $s_3$ | 0.8 | 0.09 | 2.8 | 0.252 | | 2.548 |
| $a_2$ 0.6: $s_1$ | 0.3 | 0.08 | 1.9 | 0.152 | | |
| $s_2$ | 1.0 | 0.01 | 2.6 | 0.026 | 0.202 | |
| $s_3$ | 0.8 | 0.01 | 2.4 | 0.024 | | |
| $s_3$ (0.1) $a_1$ 1.0: $s_1$ | 0.3 | 0.08 | 2.3 | 0.184 | | |
| $s_2$ | 1.0 | 0.01 | 3.0 | 0.03 | **0.242** | |
| $s_3$ | 0.8 | 0.01 | 2.8 | 0.028 | | |
| $a_2$ 0.6: $s_1$ | 0.3 | 0.01 | 1.9 | 0.019 | | |
| $s_2$ | 1.0 | 0.0 | 2.6 | 0 | 0.235 | |
| $s_3$ | 0.8 | 0.09 | 2.4 | 0.216 | | |

Figure 4 : all two-stage behaviors from $s_3$ and selection of maximum branch

$$\left( V^1(s_3) = \max_{u_1,u_2} \sum_{(x_2,x_3)\in X\times X} \{[r_1(u_1)+r_2(u_2)+r_G(x_3)]p(x_2|s_3,u_1)p(x_3|x_2,u_2)\} \right)$$

| history | ter. | path | sum | mult. | sub. | total |
|---|---|---|---|---|---|---|
| $s_1$ $a_1$ 1.0 → 0.8 $s_1$ | 0.3 | 0.64 | 2.0 | 1.28 | | |
| → 0.1 $s_2$ | 1.0 | 0.08 | 2.7 | 0.216 | 1.696 | |
| → 0.1 $s_3$ | 0.8 | 0.08 | 2.5 | 0.2 | | |
| $a_2$ 0.6 → 0.1 $s_1$ | 0.3 | 0.08 | 1.6 | 0.128 | | |
| → 0.9 $s_2$ | 1.0 | 0.72 | 2.3 | 1.656 | **1.784** | |
| → 0.0 $s_3$ | 0.8 | 0.0 | 2.1 | 0 | | |
| $s_2$ $a_1$ 1.0 → 0.0 $s_1$ | 0.3 | 0.0 | 2.0 | 0 | | 2.248 |
| → 0.1 $s_2$ | 1.0 | 0.01 | 2.7 | 0.027 | **0.252** | |
| → 0.9 $s_3$ | 0.8 | 0.09 | 2.5 | 0.225 | | |
| $a_2$ 0.6 → 0.8 $s_1$ | 0.3 | 0.08 | 1.6 | 0.128 | | |
| → 0.1 $s_2$ | 1.0 | 0.01 | 2.3 | 0.023 | 0.172 | |
| → 0.1 $s_3$ | 0.8 | 0.01 | 2.1 | 0.021 | | |
| $s_3$ $a_1$ 1.0 → 0.8 $s_1$ | 0.3 | 0.08 | 2.0 | 0.16 | | |
| → 0.1 $s_2$ | 1.0 | 0.01 | 2.7 | 0.027 | **0.212** | |
| → 0.1 $s_3$ | 0.8 | 0.01 | 2.5 | 0.025 | | |
| $a_2$ 0.6 → 0.1 $s_1$ | 0.3 | 0.01 | 1.6 | 0.016 | | |
| → 0.0 $s_2$ | 1.0 | 0.0 | 2.3 | 0 | 0.205 | |
| → 0.9 $s_3$ | 0.8 | 0.09 | 2.1 | 0.189 | | |
| $s_1$ $a_1$ 1.0 → 0.8 $s_1$ | 0.3 | 0.08 | 2.3 | 0.184 | | |
| → 0.1 $s_2$ | 1.0 | 0.01 | 3.0 | 0.03 | 0.242 | |
| → 0.1 $s_3$ | 0.8 | 0.01 | 2.8 | 0.028 | | |
| $a_2$ 0.6 → 0.1 $s_1$ | 0.3 | 0.01 | 1.9 | 0.019 | | |
| → 0.9 $s_2$ | 1.0 | 0.09 | 2.6 | 0.234 | **0.253** | |
| → 0.0 $s_3$ | 0.8 | 0.0 | 2.4 | 0 | | |
| $s_2$ $a_1$ 1.0 → 0.0 $s_1$ | 0.3 | 0.0 | 2.3 | 0 | | 2.431 |
| → 0.1 $s_2$ | 1.0 | 0.0 | 3.0 | 0 | **0** | |
| → 0.9 $s_3$ | 0.8 | 0.0 | 2.8 | 0 | | |
| $a_2$ 0.6 → 0.8 $s_1$ | 0.3 | 0.0 | 1.9 | 0 | | |
| → 0.1 $s_2$ | 1.0 | 0.0 | 2.6 | 0 | 0 | |
| → 0.1 $s_3$ | 0.8 | 0.0 | 2.4 | 0 | | |
| $s_3$ $a_1$ 1.0 → 0.8 $s_1$ | 0.3 | 0.72 | 2.3 | 1.656 | | |
| → 0.1 $s_2$ | 1.0 | 0.09 | 3.0 | 0.27 | **2.178** | |
| → 0.1 $s_3$ | 0.8 | 0.09 | 2.8 | 0.252 | | |
| $a_2$ 0.6 → 0.1 $s_1$ | 0.3 | 0.09 | 1.9 | 0.171 | | |
| → 0.0 $s_2$ | 1.0 | 0.0 | 2.6 | 0 | 2.115 | |
| → 0.9 $s_3$ | 0.8 | 0.81 | 2.4 | 1.944 | | |

In Figures 2, 3 and 4 we use the following notations:

history $= x_1$   $r_1(u_1)$ / $u_1$   $p(x_2 \mid x_1, u_1)$   $x_2$   $r_2(u_2)$ / $u_2$   $p(x_3 \mid x_2, u_2)$   $x_3$

ter. = terminal value $= r_G(x_3)$

path = path probability $= p(x_2 \mid x_1, u_1)p(x_3 \mid x_2, u_2)$

sum = sum of the three $= r_1(u_1) + r_2(u_2) + r_G(x_3)$

mult. = path × sum

sub. = subtotal expected value

total = total expected value.

## Acknowledgments

## References

Bellman, R.E. and Zadeh, L.A. (1970), Decision-making in a fuzzy environment, *Management Sci.*, 17:B141–B164.

Bertsekas, D.P. and Shreve, S.E. (1978), *Stochastic Optimal Control*, Academic Press, New York.

Howard, R.A. (1960), *Dynamic Programming and Markov Processes*, MIT Press, Cambridge, Mass.

Iwamoto, S. and Fujita, T. (1995), Stochastic decision-making in a fuzzy environment, *J. Operations Res. Soc. Japan*, 38:467–482.

Iwamoto, S., Tsurusaki, K. and Fujita, T. , On Markov policies for minimax decision processes, *submitted.*

Puterman, M.L. (1994), *Markov Decision Processes : discrete stochastic dynamic programming*, Wiley & Sons, New York.

Stokey, N.L. and Lucas Jr., R.E. (1989), *Recursive Methods in Economic Dynamics*, Harvard Univ. Press, Cambridge, MA.