J. Operations Research Soc. of Japan Vol. 18, No. 3 & No. 4, September 1975 © 1975 The Operations Research Society of Japan

# LINEAR PROGRAMMING ON RECURSIVE ADDITIVE DYNAMIC PROGRAMMING

#### SEIICHI IWAMOTO

Kyushu University

(Received October 26, 1973; Revised March 19, 1974)

### Abstract

We study, by using linear programming (LP), an infinitehorizon stochastic dynamic programming (DP) problem with the recursive additive reward system. Since this DP problem has discount factors which may depend on the transition, it includes the "discounted" Markovian decision problem. It is shown that this problem can also be formulated as one of LP problems and that the optimal stationary policy can be obtained by the usual LP method. Some interesting examples of DP models and their mumerical solutions by LP algorithm are illustrated. Furthermore, it is verified that these solutions coincides with ones obtained by Howard's policy iteration algorithm.

### 1. Introduction

We are concerned a certain class of the discrete, stochastic and infinite-horizon  $DP'_{S}$ . In general DP problems, the word "reward" or "return" is to be understood in a very broad sense ; it is not limited to any particular economic connotation (see [1 ; pp.74]). In some cases, for example, in the fields of engineering we shall be concerned the maximizing some sort of summation of reward [2 ; pp.58, 59, 102]. From this view point, Nemhauser [9 ; Chap.II-IV] introduced the deterministic  $DP'_{S}$  with recursive (not necessarily additive) return. In this paper we use the "reward system" (RS) in stead of the "return". He also treated the stochastic  $DP'_{S}$ . But their RS is restricted to only additive or multiplicative one [9 ; pp. 152-158]. Furukawa and Iwamoto [5] have extended the continuous stochastic  $DP'_{S}$  into ones with recursive (including additive and multiplicative) RS.

In 1960, Howard [6] established the policy iteration algorithm (PIA) for the discrete stochastic DP with the discounted additive RS. Recently, the author [7] proved that Howard's PIA remains valid for the discrete stochastic DP with the recursive additive (including the discounted additive but being included by recursive) RS. This DP is a discrete, stochastic and infinite-horizon version of examples in [2; pp.58, 59, 102].

On the other hand, Manne [8] originated an approach to Markovian decision problems by LP method. Since then, LP approach has been used in order to find optimal policies for discounted Markovian, average Markovian or semi-Markovian decision problems by D'Epenoux [4], De Ghellinck and Eppen [3] and Osaki and Mine [10, 11].

In this paper we shall discuss DP with recursive additive RS (hereafter abbreviated as "recursive additive DP") by LP method. In section 2, we describe this DP and give some preliminary notations and definitions used throughout this paper. In section 3, we give a formulation of this DP problem into a LP problem and show a correspondence between solutions of two problems. Section 4 is devoted to illustrate numerical examples by LP. It is shown that the optimal solution by LP algorithm is the same as one by the algorithm in [7]. Further comments are given in the last section. The method used in our proofs of results is mainly due to that of [3].

# 2. Description of recursive additive DP

A recursive additive DP is defined by six-tuple {S, A, p, r,  $\beta$ , t}. S = {1, 2, ..., N} is a set of states, A = (A<sub>1</sub>, A<sub>2</sub>, ..., A<sub>N</sub>) is an N-tuple, each A<sub>i</sub> = {1, 2, ..., K<sub>i</sub>} is a set of actions available at state i  $\in$  S, p = (p<sup>k</sup><sub>1</sub>) is a

### *Copyright* © *by ORSJ. Unauthorized reproduction of this article is prohibited.*

transition law, that is,

$$\sum_{j=1}^{N} p_{ij}^{k} = 1, p_{ij}^{k} \ge 0, \qquad i \in S, j \in S, k \in A_{i},$$

 $r = (r_{ij}^{k}), i, j \in S, k \in A_{i} \text{ is a set of stage-wise reward,}$   $\beta = (\beta_{ij}^{k}), i, j \in S, k \in A_{i} \text{ is a generalized accumulator whose}$   $\text{element } \beta_{ij}^{k} \text{ is a discount factor depending on transition}$   $(i, k, j), \text{ and } t \text{ is a translator from } R^{1} \text{ to } R^{1}.$ 

Throughout this paper we call the recursive additive DP defined by {S, A, p, r,  $\beta$ , t} simply "recursive additive DP". We sometimes use the convenient notations  $\beta(i, k, j)$ , r(i, k, j)and p(i, k, j) in stead of  $\beta_{ij}^k$ ,  $r_{ij}^k$  and  $p_{ij}^k$  respectively.

When the system starts from an initial state  $s_1 \in S$ at the l-st stage and the decision maker takes an action  $a_1 \in A_{s_1}$  on this state  $s_1$ , the system moves to the next state  $s_2 \in S$  with probability  $p(s_1, a_1, s_2)$  at the 2-nd stage and it yields a stage-wise reward  $r(s_1, a_1, s_2)$  and a discount factor  $\beta(s_1, a_1, s_2)$ . However, at the end of the 1-st stage the decision maker obtains the translated reward  $t(r(s_1, a_1, s_2))$ . The system is then repeated from the new state  $s_2 \in S$  at the 2-nd stage. If he chooses an action  $a_2 \in A_{s_2}$  on state  $s_2$ , it moves to state  $s_3$  with probability  $p(s_2, a_2, s_3)$  at the 3-rd stage. Then the system also yields a stage-wise reward  $r(s_2, a_2, s_3)$  and a discount factor  $\beta(s_2, a_2, s_3)$  at the end

of the 2-nd stage and he really receives the discounted reward  $\beta(s_1, a_1, s_2) \cdot t(r(s_2, a_2, s_3))$ . Similarly at the end of the 3-rd stage he gets a reward  $\beta(s_1, a_1, s_2) \beta(s_2, a_2, s_3)$ ,  $t(r(s_3, a_3, s_4))$ . In general when he undergoes the history  $(s_1, a_1, s_2, a_2, \dots, s_n, a_n, s_{n+1})$  of the system up to the n-th stage, he is to receive a reward  $\beta(s_1, a_1, s_2) \beta(s_2, a_2, s_3) \cdots \beta(s_{n-1}, a_{n-1}, s_n) t(r(s_n, a_n, s_{n+1}))$  at the end of the n-th stage.

Furthermore, the process goes on the (n+1)-st stage, the (n+2)-nd stage and so on.

Since we are considering a sequential nonterminating decision process, the decision maker continues to take actions infinitely. Consequently if he undergoes the history  $h = (s_1, a_1, s_2, a_2, ...)$ , he is to receive the recursive additive reward

$$\begin{aligned} v(h) &= t(r(s_1, a_1, s_2)) + \beta(s_1, a_1, s_2)t(r(s_2, a_2, s_3)) \\ &+ \beta(s_1, a_1, s_2) \beta(s_2, a_2, s_3)t(r(s_3, a_3, s_4)) \\ &+ \cdots + \beta(s_1, a_1, s_2) \beta(s_2, a_2, s_3) \cdots \beta(s_{n-1}, s_{n-1}, s_n)t(r(s_n, a_n, s_{n+1})) + \cdots \end{aligned}$$

We call V = V(h) recursive additive RS([7]). The decision maker wishes to maximize his expected reward

over the infinite future.

We are assumed that he has a complete information on his history consisted of states and actions up to date and that he knows not only the stage-wise reward  $r = (r_{ij}^k)$ , its translator  $t = t(\cdot)$  and the generalized accumulator  $\beta = (\beta_{ij}^k)$  but also the recursive additivity of RS.

Let for integer  $m \ge 1$   $\Delta_m = \{(p_1, p_2, \dots, p_m);$   $\sum_{i=1}^{m} p_i = 1, p_1 \ge 0, p_2 \ge 0, \dots, p_m \ge 0\}$ . We say a sequence  $\pi = \{f_1, f_2, \dots\}$  randomized policy if  $f_n(i) \in \Delta_{K_1}$  for all  $i \in S$ ,  $n \ge 1$ . Then we write  $f_n(i)$  as a stochastic vector  $f_n(i) = (f_n^1(i), f_n^2(i), \dots, f_n^{i_1}(i))$  for  $i \in S$ ,  $n \ge 1$ . Using randomized policy  $\pi = \{f_1, f_2, \dots\}$  means that the decision maker chooses action  $k \in A_1$  with probability  $f_n^k(i)$ in state  $i \in S$  at n-th stage. A stationary randomized policy (S-randomized policy) is the randomized policy  $\pi = \{f_1, f_2, \dots\}$  such that  $f_1 = f_2 = \dots = f$ . Such a S-randomized policy is denoted by  $\pi = f^{(\infty)}$ . The randomized policy  $\pi = \{f_1, f_2, \dots\}$  is called nonrandomized if for each  $n \ge 1$  and  $i \in S$   $f_n(i)$  is degenerate at some

 $k \in A_i$ , that is,  $f_n(i) = (0, 0, \dots, 1, 0, \dots, 0)$ . We associate with each f such that  $f(i) = (f^1(i), K$ .

 $f^{2}(i), \ldots, f^{K_{i}}(i)) \in \Delta_{K_{i}}$  for  $i \in S$  (i) the NXl column vector  $\overline{r}(f)$  whose i-th element  $\overline{r}(f)(i)$  is

$$\overline{r}(f)(i) = \sum_{k \in A_i} \sum_{j \in S} p^k_{ij} t(r^k_{ij}) f^k(i), \quad i \in S,$$

and (ii) the NXN matrix  $\overline{P}(f)$  whose (i,j) element  $\overline{P}(i,j)$  is

$$\overline{P}(f)(i,j) = \sum_{k \in A_i} p_{ij}^k p_{ij}^k f^k(i), \qquad i,j \in S.$$

If the decision maker uses a randomized policy  $\mathcal{R} = \{f_1, f_2, ...\}$ and the system starts in  $i \in S$  at 1-st stage, his recursive additive expected reward from  $\mathcal{R}$  is the column vector

$$\mathbb{V}(\mathbf{T}) = \sum_{n=0}^{\infty} \overline{\mathbb{P}}_{n}(\mathbf{T}) \overline{\mathbf{r}}(\mathbf{f}_{n+1}),$$

where  $\overline{P}_{O}(\mathbf{R}) = \mathbf{I}$ , the N  $\times$  N identity matrix, and for  $n \ge 1$ 

$$\overline{\mathbb{P}}_{n}(\pi) = \overline{\mathbb{P}}(\mathfrak{f}_{1})\overline{\mathbb{P}}(\mathfrak{f}_{2}) \cdots \overline{\mathbb{P}}(\mathfrak{f}_{n}).$$

That is, i-th element of  $V(\eta)$  is

$$\mathbb{V}(\mathbf{R})(\mathbf{i}) = \overline{\mathbf{r}}(\mathbf{f}_1)(\mathbf{i}) + \sum_{\mathbf{k} \in A_1, \mathbf{j} \in \mathbb{S}} p_{\mathbf{i}\mathbf{j}}^{\mathbf{k}} p_{\mathbf{i}\mathbf{j}}^{\mathbf{k}} f_1^{\mathbf{k}}(\mathbf{i}) \overline{\mathbf{r}}(\mathbf{f}_2)(\mathbf{j})$$

$$+\underbrace{\sum_{k \in A_{j}, j \in S, m \in A_{j}, l \in S} p_{ij}^{k} p_{jl}^{m} p_{ij}^{k} p_{jl}^{m} f_{l}^{k}(i) f_{2}^{m}(j) \overline{r}(f_{3})(l) +}$$

$$\cdots + \underbrace{\frac{p_{k \in A_{j}, j \in S, m \in A_{j}, l \in S, \ldots, t \in A_{r}, s \in S}}_{k \in A_{j}, j \in S, m \in A_{j}, l \in S, \ldots, t \in A_{r}, s \in S} p_{ij}^{k} p_{jl}^{m} \cdots p_{rs}^{t} p_{ij}^{k} p_{jl}^{m}$$

$$\dots \mathfrak{f}_{rs}^{t} \mathfrak{f}_{1}^{k}(\mathbf{i}) \mathfrak{f}_{2}^{m}(\mathbf{j}) \cdots \mathfrak{f}_{n}^{t}(\mathbf{r}) \overline{\mathfrak{r}}(\mathfrak{f}_{n+1})(\mathbf{s}) + \cdots$$

3. Formulation and algorithm by LP  
Let {S, A, p, r, 
$$\beta$$
, t} be a fixed recursive additive

DP defined at section 2, and  $\alpha' = (\alpha_1, \alpha_2, \dots, \alpha_N)$  a fixed initial (at 1-st stage) distribution of state, that is,

$$\sum_{i=1}^{N} \alpha_{i} = 1, \quad \alpha_{i} \geq 0, \qquad i=1,2,\ldots, N.$$

Let  $\{\mu_i^k(n); n \ge 1, k \in A_i, i \in S\}$  be any set of nonnegative numbers satisfying the recursive relation;

(1) 
$$\sum_{\substack{l \in A_j}} \mu_j^l(n) = \begin{cases} & \forall j, & n=1, j \in S, \\ \\ & \sum_{i \in S} \sum_{k \in A_i} \beta_{ij}^k p_{ij}^k \mu_i^k(n-1), & n \ge 2, j \in S. \end{cases}$$

In the remainder of this paper we shall assume the following assumption :

ASSUMPTION (I).  $0 \leq \rho_{ij}^k < 1$  for any i, jeS, keA<sub>i</sub>.

LEMMA 3.1. Under the Assumption (I), any nonnegative  $\{\mu_1^k(n) ; n \ge 1, k \in A_i, i \in S\}$  satisfying (1) has the following properties :

(i) 
$$\sum_{i \in S} \sum_{k \in A_i} \mu_i^k(1) = 1, \sum_{k \in A_i} \mu_i^k \ge 0,$$
 ies,

(ii) 
$$\beta_*^{n-1} \leq \sum_{i \in S} \sum_{k \in A_i} \mu_i^k(n) \leq \beta^{*n-1}, n \geq 2.$$

Therefore, we have

Copyright © by ORSJ. Unauthorized reproduction of this article is prohibited.

$$\frac{r_{*}}{1-\beta_{*}} \leq \sum_{n \geq 1} \sum_{i \in S} \sum_{k \in A_{i}} \sum_{j \in S} p_{ij}^{k} t(r_{ij}^{k}) \mu_{i}^{k}(n) \left\langle \frac{r^{*}}{1-\beta^{*}} \right\rangle,$$

where 
$$\mathbf{r}_{*} = \min_{\substack{i,j \in S, k \in A_{i}}} t(\mathbf{r}_{ij}^{k}), \mathbf{r}^{*} = \max_{\substack{i,j \in S, k \in A_{i}}} t(\mathbf{r}_{ij}^{k})$$
  
 $\beta_{*} = \min_{\substack{i,j \in S, k \in A_{i}}} \beta_{ij}^{k}$  and  $\beta^{*} = \max_{\substack{i,j \in S, k \in A_{i}}} \beta_{ij}^{k}$ .

PROOF. Property (i) is a trivial consequence. Property (ii) is to be proved by induction on n.

LEMMA 3.2. Under the Assumption (I), (i) any randomized policy  $\eta = \{f_1, f_2, \cdots\}$  gives a nonnegative solution  $\{\mu_1^k(n)\}$  of (1) and vice versa, and, furthermore,

(ii) 
$$\sum_{i \in S} \alpha_i V(\pi)(\iota) = \sum_{n=1}^{\infty} \sum_{i \in S} \sum_{k \in A_1} \overline{r}_1^k \mu_1^k(n),$$
  
where  $\overline{r}_1^k = \sum_{j \in S} p_{ij}^k t(r_{ij}^k).$ 

PROOF. Let  $\mathcal{T} = \{f_1, f_2, \cdots\}$  be any randomized policy. Then we can give a nonnegative  $\mu_j^{\ell}(n)$  for  $n \ge 1$ ,  $\ell \in A_j$ ,  $j \in S$  as follows ;

(1)'  
$$\begin{cases} \mu_{j}^{\ell}(1) = \alpha_{j} f_{1}^{\ell}(j), & \ell \in A_{j}, j \in S, \\ \mu_{j}^{1}(n+1) = \sum_{i \in S} \sum_{k \in A_{i}} p_{ij}^{k} \rho_{ij}^{k} \mu_{i}^{k}(n) f_{n+1}^{1}(j), \quad n \ge 1, \end{cases}$$

 $l \in A_i$ ,  $j \in S$ .

Obviously, these  $\{\mu_j^1(n) ; n \ge 1, l \in A_{\epsilon}, j \in S\}$  satisfy (1).

Conversely, let nonnegative  $\{\mu_1^k(n)\}$  satisfy (1). Then, we can define f<sub>n</sub> as follows ;

$$\begin{cases} f_{1}^{k}(i) = \frac{\mu_{1}^{k}(1)}{\alpha_{1}}, & n=1, \ k \in A_{1}, \ i \in S, \\ f_{n}^{1}(j) = \frac{\mu_{j}^{1}(n)}{\sum_{i \in S} \sum_{k \in A_{i}} \beta_{i,j}^{k} p_{i,j}^{k} \mu_{1}^{k}(n-1)}, & n \ge 2, \ l \in A_{j}, \ j \in S, \end{cases}$$

where  $\frac{0}{0} = 0$ . Then the policy  $\Re = \{f_1, f_2, \dots\}$  is a randomized policy. Moreover, we have, by using (1)' and exchanging the summation,  $\sum_{i \in S} \sum_{k \in A_i} \overline{r}_i^k \mu_i^k(n+1)$ 

$$= \underbrace{\mathbb{Z}}_{i \in S} \alpha_{i} \xrightarrow{k \in A_{i}, j \in S, m \in A_{j}, l \in S, \dots, t \in A_{r}, s \in S} p_{ij}^{k} p_{jl}^{m} \cdots$$

$$p_{rs}^{t} \beta_{ij}^{k} \beta_{jl}^{m} \cdots \beta_{rs}^{t} f_{1}^{k}(i) f_{2}^{k}(j) \cdots f_{n}^{t}(r) \overline{r}(f_{n+1})(s), n \ge 0.$$

Hence (ii) holds. This completes the proof.

We note that 
$$\sum_{n=1}^{\infty} \sum_{i \in S} \sum_{k \in K_{1}} \overline{r}_{1}^{k} \mu_{1}^{k}(n)$$

is the total expected recursive additive reward obtained from the randomized policy  $\mathcal{T} = \{f_1, f_2, \cdots\}$  corresponding  $\{\mu_i^k(n)\}$ , started in the initial distribution &. Consequently, above lemmas and note enable us to give a maximization problem  $(P_0)$ :

Problem (P<sub>0</sub>) : Maximize 
$$\sum_{n=1}^{\infty} \sum_{i \in S} \sum_{k \in A_i} \overline{r}_i^k \mu_i^k(n)$$

(1) 
$$\sum_{\boldsymbol{\ell} \in \mathbf{A}_{j}}^{\text{under}} \mu_{j}^{\boldsymbol{\ell}}(n) = \begin{cases} \alpha_{j}, & n=1, \\ \\ \sum_{i \in S} \sum_{k \in \mathbf{A}_{i}} \beta_{ij}^{k} p_{ij}^{k} \mu_{i}^{k}(n-1), & n \geq 2, j \in S, \end{cases}$$

(2) 
$$\mu_{i}^{k}(n) \ge 0$$
,  $n \ge 1$ ,  $k \in A_{i}$ ,  $i \in S$ .

By Lemma 3.1, we can define a set of the new variables  $\{y_1^k\}$  as follows :

$$y_{i}^{k} = \sum_{n=1}^{\infty} \mu_{1}^{k}(n), \qquad k \in A_{i}, i \in S.$$

Hence, we have a modified maximization problem  $({\rm P}_{\rm T})$  :

Problem ( $P_{\tau \tau}$ ) : Maximize

(3) 
$$\sum_{i \in S} \sum_{k \in A_i} \overline{r}_i^{k} y_i^{k}$$

under

(4) 
$$\sum_{\boldsymbol{l} \in A_{j}} y_{j}^{\boldsymbol{l}} - \sum_{i \in S} \sum_{k \in A_{i}} \beta_{ij}^{k} p_{ij}^{k} y_{i}^{k} = \alpha_{j}, \quad j \in S,$$
  
(5) 
$$y_{i}^{k} \ge 0, \quad k \in A_{i}, \quad i \in S.$$

Next lemma states the relationship between Problem (P\_0) and Problem (P\_m).

LEMMA 3.3. If  $\{\mu_{i}^{k}(n)\}$  is a nonnegative solution of (1), then  $\{y_{i}^{k}\}$  is a solution of Problem  $(P_{T})$ , and  $\sum_{i\in S} \sum_{k\in A_{i}} \overline{r}_{i}^{k} y_{i}^{k}$  is the expected recursive additive reward

which corresponds to  $\left\{ \mu_{1}^{k}(n) 
ight\}$  .

PROOF. It is easy to show that  $\{y_1^k\}$  satisfies (4) and (5).

We can define a S-nonrandomized policy  $\pi = f^{(\infty)}$  by a

function f such that for each  $i \in S$  selects exactly one variable  $y_1^k$   $k \in A_1$ . This fact is easy to check.

THEOREM 3.1. Let Assumption (I) be satisfied. If the equation (4) is restricted to the variables  $y_1^k$  selected by any S-nonrandomized policy, then : (i) the corresponding subsystem has a unique solution,

(ii) if  $\alpha_1 \ge 0$  is then  $y_1^k \ge 0$  is, (iii) if  $\alpha_1 \ge 0$  is, then  $y_1^k \ge 0$  is.

PROOF. This theorem corresponds to Proposition 2.3 in [3] which treated the case of  $\beta_{ij}^k \equiv \beta$ . The proof is similar to that of Proposition 2.3.

LEMMA 3.4. Let Assumption (I) be satisfied and  $\alpha_1 > 0$ for i $\in$ S. Then there exists an one to one correspondence between S-nonrandomized policies and basic feasible solutions of (4), (5). Moreover, any basic feasible solution is nondegenerate.

PROOF. The proof follows in the same way as in Proposition 2.4 of [3], and the details are omitted.

Lemma 3.4 yields the following definition of optimality. A S-nonrandomized policy  $\pi = f^{(\infty)}$  is optimal if its corresponding basic feasible solution is optimal.

THEOREM 3.2. Let Assumption (I) be satisfied. Whenever  $\alpha_1 > 0$  for i $\in$ S, the Problem (P<sub>T</sub>) has an optimal basic solution and its dual problem has a unique optimal solution. Any optimal S-policy associated with it remains optimal for any ( $\alpha_1, \alpha_2, \dots, \alpha_N$ ) such that  $\alpha_1 > 0$  for i $\in$ S.

PROOF. The proof is similar to that of Proposition 3.5 of [3], and the details are omitted.

COROLLARY For  $\aleph_1 > 0$  for  $i \in S$  (say  $\aleph_1 = \frac{1}{N}$ ,  $i \in S$ ) there exists an optimal basic solution such that for each  $i \in S$  there is exactly one k such that  $y_1^k > 0$  and  $y_1^k = 0$ for k otherwise.

PROOF. This is a straightforward from Lemma 3.4 and Theorem 3.2.

#### 4. Numerical examples

We now illustrate correspondence between the optimal solution by PIA and the optimal solution by LP algorithm. As for the definition, reward system and optimal solution by PIA of the following DP3, see the corresponding example in [5].

EXAMPLE 1 (General Additive DP)

In the general additive DP { S, A, p, r,  $\beta$ }, the objective function is the expected value of the general additive RS

 $V(h) = r_1 + p_1 r_2 + \beta_1 p_2 r_3 + \dots + \beta_1 p_2 \cdots \beta_{n-1} r_n + \dots,$ 

since this is the case where t(r) = r in the recursive additive DP{S, A, p, r, p, t}. Following data is a slightly modified one from Howard [4]. Of course Assumption (I) is satisfied.

state action		tr	transition probability		 stage-wise reward			generalized accumulator		
i	k	p <sup>k</sup> il	p <sup>k</sup> 12	p <sup>k</sup> i3	r <sup>k</sup> il	r <sup>k</sup> i2	r <sup>k</sup> i3	₿ <sup>k</sup> il	/ <sup>k</sup> 12	k /13
1	1	1 2	1 Ę	1 1 1	 10	4	8	.95	.98	.98
	2	$\frac{1}{16}$	<u>3</u> 4	<u>3</u> 16	8	2	4	.90	.90	•93
	3	<u>1</u> 4	<u>1</u> 8	<u>5</u> 8	4	6	4	.98	.96	.98
2	l	$\frac{1}{2}$	0	$\frac{1}{2}$	14	0	18	.85	.90	•95
	2	$\frac{1}{16}$	<u>7</u> 8	$\frac{1}{16}$	6	16	8	.80	.80	•95
	3	$\frac{1}{3}$	<u>1</u> 3	<u>1</u> 3	<del>-</del> 5	-5	-5	•95	•95	•95
3	l	$\frac{1}{4}$	<u>1</u> 4	$\frac{1}{2}$	10	2	8	•75	.90	•95
	2	$\frac{1}{8}$	<u>3</u> 4	$\frac{1}{8}$	6	4	2	•95	.70	.80
	3	<u>3</u> 4	$\frac{1}{16}$	<u>3</u> 16	4	0	8	•95	•95	۰95

TABLE 4.1. Data for general additive DP

Then PIA yields an optimal S-policy  $f^{(\infty)}$ , where  $f = \begin{pmatrix} 1 \\ 1 \\ 3 \end{pmatrix}$  and an optimal

return  $V(f^{(\infty)}) = \begin{pmatrix} 169 & 490 \\ 166 & 129 \\ 164 & 411 \end{pmatrix}$ .

On the other hand, for an initial vector  $\alpha = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ the LP Problem (P<sub>T</sub>) becomes :

Maximize  $8y_1^1 + \frac{11}{4}y_1^2 + \frac{17}{4}y_1^3 + 16y_2^1 + \frac{31}{2}y_2^2 + (15)y_2^3 + 7y_3^1 + 4y_3^2 + \frac{9}{2}y_3^3$ 

subject to

 $\frac{105}{200} \frac{1}{1} + \frac{1510}{1600} \frac{2}{1} + \frac{302}{400} \frac{3}{1} - \frac{85}{200} \frac{1}{2} - \frac{80}{1600} \frac{2}{2} - \frac{95}{300} \frac{3}{2}$  $- \frac{75}{400} \frac{1}{3} - \frac{95}{800} \frac{2}{3} - \frac{285}{400} \frac{3}{3} = \frac{1}{3},$  $- \frac{98}{400} \frac{1}{1} - \frac{270}{400} \frac{2}{1} - \frac{96}{800} \frac{3}{1} + \frac{1}{92} + \frac{30}{100} \frac{2}{2} + \frac{205}{300} \frac{3}{2}$  $- \frac{90}{400} \frac{1}{3} - \frac{210}{400} \frac{2}{3} - \frac{95}{1600} \frac{3}{3} = \frac{1}{3},$  $- \frac{98}{400} \frac{1}{1} - \frac{279}{400} \frac{2}{1} - \frac{96}{800} \frac{3}{1} + \frac{210}{400} \frac{2}{3} - \frac{95}{1600} \frac{3}{3} = \frac{1}{3},$  $- \frac{98}{400} \frac{1}{1} - \frac{279}{1600} \frac{2}{1} - \frac{490}{800} \frac{3}{1} - \frac{95}{200} \frac{2}{2} - \frac{95}{300} \frac{3}{2} = \frac{95}{300} \frac{3}{2}$ 

 $+\frac{105}{200}y_3^1 + \frac{90}{100}y_3^2 + \frac{1315}{1600}y_3^3 = \frac{1}{3},$ 

 $y_1^1, y_1^2, y_1^3, y_2^1, y_2^2, y_2^3, y_3^1, y_3^2, y_3^3 \ge 0.$ 

The optimal solution of this LP problem is

$$(y_1^1, y_1^2, y_1^3, y_2^1, y_2^2, y_2^3, y_3^1, y_3^2, y_3^3)$$

= (10.9688, 0.0, 0.0, 3.3540, 0.0, 0.0, 0.0, 0.0, 5.6138)

and its (optimal) value of the objective function is 166.6768. Note that this value is nearly equal to

Furthermore this optimal solution shows that  $f = \begin{pmatrix} 1 \\ 1 \\ 3 \end{pmatrix}$  is optimal.

EXAMPLE 2(Multiplicative Additive DP)

The multiplicative additive DP { S, A, p, r } is the case where  $\beta_{ij}^{k} \equiv r_{ij}^{k}$ , t(r) = r in the recursive additive DP. Then, the objective function of this DP { S, A, p, r } is the expected value of the multiplicative additive RS

 $V(h) = r_1 = r_1r_2 + r_1r_2r_3 + \dots + r_1r_2 \dots + r_n + \dots$ 

The following data satisfies Assumption (I).

TABLE 4.2.

Data for multiplicative additive DP

state	action	trans pr	itio obab	n ility		sta r	te-wi eward	se	-
i	k	p <sup>k</sup> il	p <sup>k</sup> 12	p <sup>k</sup> 13	-	r <sup>k</sup> il	r <sup>k</sup> 12	r <sup>k</sup> i3	
l	l	<u>1</u> 2	1 4	1 4		$\frac{1}{2}$	<u>1</u> 5	2 5	
	2	1 16	<u>3</u> 4	<u>3</u> 16		<u>2</u> 5	$\frac{1}{10}$	<u>1</u> 5	
2	l	<u>1</u> 2	0	$\frac{1}{2}$		$\frac{7}{10}$	1 20	<u>9</u> 10	
	2	$\frac{1}{16}$	<u>7</u> 8	$\frac{1}{16}$		<u>2</u> 5	<u>4</u> 5	<u>2</u> 5	
	3	$\frac{1}{3}$	<u>1</u> 3	$\frac{1}{3}$		$\frac{1}{20}$	$\frac{1}{20}$	$\frac{1}{20}$	
3	1	1 4	1 4	12		1 2	$\frac{1}{10}$	2 5	
	2	1 8	3 4	<u>1</u> 8		$\frac{3}{10}$	<u>1</u> 5	$\frac{1}{10}$	
	3	3 4	$\frac{1}{16}$	<u>3</u> 16		<u>1</u> 5	$\frac{1}{20}$	2 5	
							·		

Then, by PIA, we have an optimal S-policy  $f^{(\infty)}$ , where  $f = \begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix}$ , and the optimal return

 $V(f^{(\infty)}) = \begin{pmatrix} 0.7938\\ 2.6198\\ 0.6434 \end{pmatrix}.$ 

The LP problem  $(P_T)$  for  $\chi = (\frac{1}{4}, \frac{1}{4}, \frac{1}{2})$  has an optimal solution  $(y_1^1, y_1^2, y_1^3, y_2^1, y_2^2, y_2^3, y_3^1, y_3^2, y_3^3) = (0.4851, 0.0, 0.0, 0.0, 0.9739, 0.0, 0.7161, 0.0, 0.0)$  and an optimal value 1.1751. Note that this optimal value is equal

EXAMPLE 3 (Divided Additive DP)

The divided additive DP $\{S, A, p, r\}$  has the divided additive RS

$$V(h) = r_1 + \frac{r_2}{r_1} + \frac{r_3}{r_1 r_2} + \cdots + \frac{r_n}{r_1 r_2 \cdots r_{n-1}} + \cdots$$

since this is the case where  $\beta_{ij}^k \equiv 1/r_{ij}^k$ , t(r) = r in the recursive additive DP. We can illustrate a DP with  $\beta_{1j}^k \equiv 1/r_{1j}^k$ ,  $r_{1j}^k \equiv k$ ,  $t(r) = r^b$  (b>0) in [2;pp.58]. This DP has con-

tinuous state-action spaces, deterministic transition law and finite horizon. In the divided additive DP Assumption (I) means  $r_{ij}^k > 1$  for ies, keA<sub>i</sub>, jes, which is satisfied by the following data.

### Copyright © by ORSJ. Unauthorized reproduction of this article is prohibited.

state	action	transition probability	stage-wise reward
i	k	$p_{i1}^k p_{i2}^k p_{i3}^k$	r <sup>k</sup> r <sup>k</sup> r <sup>k</sup> r <sup>k</sup> il i2 i3
l	l	$\frac{1}{2}$ $\frac{1}{4}$ $\frac{1}{4}$	$\frac{3}{2}$ $\frac{6}{5}$ $\frac{7}{5}$
	2	$\frac{1}{16}$ $\frac{3}{4}$ $\frac{3}{16}$	$\frac{7}{5}$ $\frac{11}{10}$ $\frac{6}{5}$
	3	1 1 5 4 8 8	$\frac{6}{5}$ $\frac{13}{10}$ $\frac{6}{5}$
2	l	$\frac{1}{2}$ 0 $\frac{1}{2}$	$\frac{17}{10}$ $\frac{21}{20}$ $\frac{19}{10}$
	2	$\frac{1}{16}$ $\frac{7}{8}$ $\frac{1}{16}$	<u>7 9 26</u> 5 5 25
	3	$\frac{1}{3}$ $\frac{1}{3}$ $\frac{1}{3}$	$\frac{21}{20}$ $\frac{21}{20}$ $\frac{21}{20}$
3	1	$\frac{1}{4}$ $\frac{1}{4}$ $\frac{1}{2}$	$\frac{3}{2}$ $\frac{11}{10}$ $\frac{7}{5}$
	2	1 <u>3</u> 1 8 <u>4</u> 8	$\frac{13}{10}$ $\frac{6}{5}$ $\frac{11}{10}$
	3	$\frac{3}{4}$ $\frac{1}{16}$ $\frac{3}{16}$	$\frac{6}{5}$ $\frac{21}{20}$ $\frac{7}{5}$

TABLE 4.3. Data for divided additive DP

Then optimal S-policy is specified by  $f = \begin{pmatrix} 2\\3\\2 \end{pmatrix}$ , and optimal return is  $V(f^{(\infty)}) = \begin{pmatrix} 11.8020\\12.2804\\11.2934 \end{pmatrix}$ .

If process strates at initial distribution  $\chi = (\frac{1}{5}, \frac{2}{5}, \frac{2}{5})$ , the LP algorithm yields an optimal solution  $(y_1^1, y_1^2, y_1^3, y_2^1, y_2^2, y_2^3, y_3^1, y_3^2, y_3^3) = (0.0, 2.3176, 0.0, 0.0, 0.0, 5.4885, 0.0, 2.8256, 0.0)$  and an optimal value 11.7899. Note that

EXAMPLE 4 (Exponential Additive DP)

The exponential additive DP {S, A, p, r} has the exponential additive RS

$$V(h) = r_1 + e^{r_1} \cdot r_2 + e^{r_1 + r_2} \cdot r_3 + \dots + e^{r_1 + r_2 + \dots + r_{n-1}} \cdot r_n + \dots,$$

since this is the case where  $\beta_{1j}^{k} \equiv e^{r_{1j}^{k}}$ , t(r)=r in<sub>k</sub> the recursive additive DP. We have a DP with  $\beta_{1j}^{k} \equiv e^{r_{1j}}$ ,  $t(r) = (1-r)e^{r}$  [2; pp.102]. But this DP has continuous action space, deterministic transition law and finite horizon. If  $r_{1j}^{k} < 0$  for  $i \in S$ ,  $k \in A_{1}$ ,  $j \in S$  then the exponential additive DP satisfies Assumption (I). The following data satisfies Assumption (I).

TABLE 4.4.

Data for exponential additive DP

state .	action	transition probability	stage-wise reward
i	k	$p_{11}^k p_{12}^k p_{13}^k$	r <sup>k</sup> r <sup>k</sup> r <sup>k</sup> il i2 i3
1	l	$\frac{1}{2}$ $\frac{1}{4}$ $\frac{1}{4}$	$\frac{1}{2}$ $\frac{1}{5}$ $\frac{2}{5}$
	2	$\frac{1}{16}$ $\frac{3}{4}$ $\frac{3}{16}$	$\frac{2}{5}$ $\frac{1}{10}$ $\frac{1}{5}$
	3	1 1 5 4 8 8	$\frac{1}{5}$ $\frac{3}{10}$ $\frac{1}{5}$
2	1	$\frac{1}{2}$ 0 $\frac{1}{2}$	$-\frac{7}{10}$ $-\frac{1}{20}$ $-\frac{9}{10}$
	2	$\frac{1}{16}  \frac{7}{8}  \frac{1}{16}$	2 4 2
	3	$\frac{1}{3}$ $\frac{1}{3}$ $\frac{1}{3}$	$\frac{1}{20}$ $\frac{1}{20}$ $\frac{1}{20}$
3	l	$\frac{1}{4}$ $\frac{1}{4}$ $\frac{1}{2}$	$\frac{1}{2}$ $\frac{1}{10}$ $\frac{2}{5}$
	2	$\frac{1}{8}$ $\frac{3}{4}$ $\frac{1}{8}$	$\frac{3}{10} \frac{1}{5} \frac{1}{10}$
	3	$\frac{3}{4}$ $\frac{1}{16}$ $\frac{3}{16}$	$\frac{1}{5}$ $\frac{1}{20}$ $\frac{2}{5}$

We have optimal stationary policy  $f^{(\infty)}$ , where  $f = \begin{pmatrix} 2\\3\\2 \end{pmatrix}$ and optimal return  $V(f^{(\infty)}) = \begin{pmatrix} -1.0831\\-1.0807\\-1.0867 \end{pmatrix}$ .

If  $\alpha = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ , then the LP problem  $(P_{T})$  yields an optimal solution  $(y_{1}^{1}, y_{1}^{2}, y_{1}^{3}, y_{2}^{1}, y_{2}^{2}, y_{2}^{3}, y_{3}^{1}, y_{3}^{2}, y_{3}^{3}) = (0.0, 2.2768, 0.0, 0.0, 0.0, 5.0739, 0.0, 2.5839, 0.0)$  and an optimal value -1.0835. We can verify that

$$( V(f^{(\infty)}) = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}) \begin{pmatrix} -1.0831\\ -1.0807\\ -1,0867 \end{pmatrix}$$

coincides with optimal value.

EXAMPLE 5 (Logarithmic Additive DP)

This is the case where  $\beta_{ij}^k \equiv \log r_{ij}^k$ , t(r) = r in the recursive additive DP { S, A, p, r,  $\beta$ , t }. Then, the logarithmic additive RS is given as follows :

 $V(h) = r_1 + (\log r_1) \cdot r_2 + (\log r_1 \cdot \log r_2) r_3 + \cdots +$ 

 $(\log r_1 \cdot \log r_2 \cdots \log r_{n-1})r_n + \cdots$ 

In this DP Assumption (I) means that  $1 \le r_{ij}^k < e$  for  $i \in S$ ,  $k \in A_i$ ,  $j \in S$ . The following data satisfies Assumption (I).

state	action	transition probability	stage-wise reward			
i	k	$p_{i1}^k p_{i2}^k p_{i3}^k$	r <sup>k</sup> r <sup>k</sup> r <sup>k</sup> il <sup>i</sup> 2 r <sup>k</sup> i3			
1	1	$\frac{1}{2}  \frac{1}{4}  \frac{1}{4}$	2.3 2.7 2.4			
	2	$\frac{1}{16}$ $\frac{3}{4}$ $\frac{3}{16}$	2.7 2.3 2.6			
	3	1 <u>1 5</u> 4 8 8	2.5 2.4 2.6			
2	1	$\frac{1}{2}$ 0 $\frac{1}{2}$	2.7 2.3 2.4			
	2	$\frac{1}{16}  \frac{7}{8}  \frac{1}{16}$	2.6 2.4 2.7			
	3	$\frac{1}{3}$ $\frac{1}{3}$ $\frac{1}{3}$	2.4 2.1 2.5			
3	1	$\frac{1}{4}$ $\frac{1}{4}$ $\frac{1}{2}$	2.6 2.5 2.7			
	2	1 <u>31</u> 8 <u>4</u> 8	2.7 2.6 2.4			
	3	$\frac{3}{4}$ $\frac{1}{16}$ $\frac{3}{16}$	2.6 2.7 2.5			

TABLE 4.5. Data for logarithmic additive DP

Then optimal S-policy is  $f^{(\infty)}$  and optimal return is

$$V(f^{(\infty)}) = \begin{pmatrix} 52.3188\\ 52.0526\\ 53.7307 \end{pmatrix}$$
, where  $f = \begin{pmatrix} 3\\ 1\\ 1 \end{pmatrix}$ .

The LP problem  $(P_T)$  with an initial distribution  $\alpha = (\frac{1}{2}, \frac{1}{4}, \frac{1}{4})$  gives an optimal solution  $(y_1^1, y_1^2, y_1^3, y_2^1, y_2^2, y_2^3, y_3^1, y_3^2, y_3^3) = (0.0, 0.0, 6.1654, 3.3892, 0.0, 0.0, 10.7585, 0.0, 0.0)$ and an optimal value 52.6052. Note that this value is

We remark that above five examples are the case t(r)=rin the recursive additive DP { S, A, p, r,  $\beta$ , t }. But we can treat, for example, the case where  $t(r)=\frac{1}{r}$ ,  $t(r)=e^{r}$ ,  $t(r)=(1-r)e^{r}$ ,  $t(r)=\log r$ , etc., ([7]).

# 5. Further remarks

In this section we shall give some remarks on the recursive additive DP.

Let { S, A, p, r,  $\beta$ , t } be the recursive additive DP satisfying Assumption (I). We define DP {  $\overline{S}$ ,  $\overline{A}$ ,  $\overline{p}$ ,  $\overline{r}$  } in which

 $\overline{S} = SV\{0\}$ ,  $O(\ S)$  is a fictitious state,

Copyright © by ORSJ. Unauthorized reproduction of this article is prohibited.

$$\overline{A} = (A_0, A_1, \dots, A_N), A_0 = \{1\},$$

$$\vec{p}_{ij}^{k} = \begin{cases} 1, & i=0, k=1, j=0, \\ 1 - \sum_{j \in S} \beta_{ij}^{k} p_{ij}^{k}, & i\in S, k\in A_{i}, j=0, \\ \\ \beta_{ij}^{k} p_{ij}^{k}, & i\in S, k\in A_{i}, j\in S, \\ \end{cases}$$

and

$$\overline{r}_{ij}^{k} = \begin{cases} 0, & i=0, k=1, j=0, \\ t(r_{ij}^{k}), & i\in S, k\in A_{i}, j\in S. \end{cases}$$

Note that  $\overline{P}(X_{n+1}=j|X_n=i, Y_n=k)=\overline{p}_{ij}^k$  for  $i\in\overline{S}$ ,  $k\in A_i$ ,  $j\in\overline{S}$ ,

where  $\overline{P}$  is a probability law associated with  $DP\{\overline{S}, \overline{A}, \overline{p}, \overline{r}\}$ , and  $X_n, Y_n(n\geq 1)$  denote observed state and action at n-th stage. In other words, nonnegative  $\mu_1^k(n)$  satisfying (1) is the joint probability of being in state  $i \in \overline{S}$  and making decision  $k \in A_i$  at the n-th stage regarding to above probability law  $\overline{P}$ .

Furthermore above  $\{\overline{S}, \overline{A}, \overline{p}, \overline{r}\}$  gives DP with an absorved state  $\{0\}$ . We can also apply the LP method for DP $\{\overline{S}, \overline{A}, \overline{p}, \overline{r}\}$  as well as DP $\{S, A, p, r, \beta, t\}$  with Assumption (I). But it is rather difficut to get five examples in section 4 from the reduced DP $\{\overline{S}, \overline{A}, \overline{p}, \overline{r}\}$ .

Acknowledgement

The author wishes to express his hearty thanks to Prof. N. Furukawa for his advices. He also thanks the referee for his various comments and suggestions for improving this paper.

#### References

- Aris, R., Discrete Dynamic Programming, Blaisdell, Publishing Company, New York Tront London, (1964).
- [2] Bellman, R., Dynamic Programming, Princeton Univ. Press, Princeton, New Jersey, (1957).
- [3] DeGhellinck, G.T. and Eppen, G.D., "Linear programming solutions for separable Markovian decision problems", Mangt. Sci., 13, 371-394, (1967).
- [4] D'Epenoux, F., "A probabilistic production and inventory problem", Mangt. Sci., 10, 98-108 (1963).
- [5] Furukawa, N. and Iwamoto, S., "Markovian decision processes with recursive reward functions", Bull. Math. Statist., 15, 3-4, 79-91, (1973).
- [6] Howard, R.A., Dynamic Programming and Markov Processes,M.I.T. Press, Cambridge, Massachusetts, (1960).
- [7] Iwamoto, S., "Discrete dynamic programming with recursive additive system", Bull. Math. Statist., 16, 1-2, 49-66, (1974).
- [8] Manne, A.S., "Linear programming and sequential decisions", Mangt. Sci., 6, 259-267, (1960).

*Copyright* © *by ORSJ. Unauthorized reproduction of this article is prohibited.* 

- [9] Nemhauser, G.L., Introduction to Dynamic Programming, John Wiley and Sons, NewwYork London Sydney, (1966).
- [10] Osaki, S. and Mine, H., "Linear programming algorithm for semi-Markovian decision processes", J. Math. Anal. Appl., 22, 356-381, (1968).
- [11] Osaki, S. and Mine, H., "Some remarks on a Markovian decision problem with an absorbing state", J. Math. Anal. Appl., 23, 327-333, (1968).