

分数型評価のマルコフ決定過程

九大・経済 岩本 誠一
九工大・工 藤田 敏治

Abstract

We consider how to optimize a ratio of two expected values of additive statistics on a finite-state controlled Markov chain. We present an algorithm for finding an optimal policy by use of both stochastic dynamic programming and fractional programming.

1 Introduction

We are concerned with finding an optimal policy which maximizes a ratio of two expected values of additive rewards over a controlled Markov decision process ([7],[8]).

2 Fractional Expectation Problem

Throughout the paper, the following data is given:

- $N \geq 1$ is an integer; the *total number of stages*
- $S = \{s_1, s_2, \dots, s_p\}$ is a *finite state space*
- $A = \{a_1, a_2, \dots, a_k\}$ is a *finite action space*
- $r : S \times A \rightarrow R^1, R : S \times A \rightarrow (0, \infty)$ are two *n-th reward functions*
- $k : S \rightarrow R^1, K : S \rightarrow (0, \infty)$ are two *terminal reward functions* (1)
- β is a *discount factor* : $0 < \beta < 1$
- p is a *Markov transition law*
- : $p(y|x, u) \geq 0 \quad \forall (x, u, y) \in S \times A \times S, \quad \sum_{y \in S} p(y|x, u) = 1 \quad \forall (x, u) \in S \times A$
- $y \sim p(\cdot |x, u)$ denotes that next state y conditioned on state x and action u appears with probability $p(y|x, u)$.

We use the following simple notations:

$$\begin{aligned}
 r_n &:= r(X_n, U_n), \quad R_n := R(X_n, U_n) \quad 1 \leq n \leq N \\
 r_{N+1} &:= k(X_{N+1}), \quad R_{N+1} := K(X_{N+1}) \\
 E_{x_n}[Y] &:= E[Y|X_n = x_n].
 \end{aligned}
 \tag{2}$$

Let $c \in R^1$ be a given constant (level). Then we consider how to maximize the ratio of the expected value of one *additive* statistics

$$r(X_1, U_1) + r(X_2, U_2) + \dots + r(X_N, U_N) + k(X_{N+1})$$

to that of the other

$$R(X_1, U_1) + R(X_2, U_2) + \dots + R(X_N, U_N) + K(X_{N+1}).$$

A Markov policy $\pi = \{\pi_1, \pi_2, \dots, \pi_N\}$ is a finite sequence of decision functions:

$$\pi_n : S \rightarrow A \quad 1 \leq n \leq N. \quad (3)$$

The set of all Markov policies is denoted by Π . Given an initial state $x_1 \in S$, let us consider the maximization problem:

$$F(x_1) \quad \text{Maximize} \quad \frac{E_{x_1}^{\pi} \left[\sum_{n=1}^{N+1} r_n \right]}{E_{x_1}^{\pi} \left[\sum_{n=1}^{N+1} R_n \right]} \quad \text{subject to} \quad (i) \quad \pi \in \Pi. \quad (4)$$

By introducing the Lagrange multiplier λ , the fractional optimization problem (4) is transformed into the standard stochastic optimization problem with the following additive criteria:

$$\begin{aligned} P(x_1; \lambda) \quad & \text{Maximize} \quad E_{x_1}^{\pi} \left[\sum_{n=1}^{N+1} (r_n - \lambda R_n) \right] \\ & \text{subject to} \quad (i) \quad x_{n+1} \sim p(\cdot | x_n, u_n) \quad 1 \leq n \leq N \\ & \quad \quad \quad (ii) \quad u_n \in A \quad 1 \leq n \leq N \\ & \quad \quad \quad x_1 \in S, \quad \lambda \in R^1, \quad 1 \leq n \leq N+1. \end{aligned} \quad (5)$$

Let $u_n(x_n; \lambda)$ be the maximum value of the subproblem:

$$\begin{aligned} P_n(x_n; \lambda) \quad & \text{Maximize} \quad E_{x_n}^{\pi} \left[\sum_{m=n}^{N+1} (r_m - \lambda R_m) \right] \\ & \text{subject to} \quad (i) \quad x_{m+1} \sim p(\cdot | x_m, u_m) \quad n \leq m \leq N \\ & \quad \quad \quad (ii) \quad u_m \in A \quad n \leq m \leq N \\ & \quad \quad \quad x_n \in S, \quad \lambda \in R^1, \quad 1 \leq n \leq N+1. \end{aligned} \quad (6)$$

Then we have the recursive equation([4]):

THEOREM 2.1

$$\begin{aligned} u_n(x; \lambda) &= \text{Max}_{u \in A} [r(x, u) - \lambda R(x, u) + \sum_{y \in S} u_{n+1}(y; \lambda) p(y|x, u)] \\ & \quad \quad \quad x \in S, \quad \lambda \in R^1, \quad 1 \leq n \leq N \\ u_{N+1}(x; \lambda) &= k(x) - \lambda K(x) \quad x \in S, \quad \lambda \in R^1. \end{aligned} \quad (7)$$

3 Infinite-stage Problem

In this section we consider an optimization problem of the ratio of one total discounted expected value over an infinite-stage to the other as follows:

$$F'(x_1) \quad \text{Maximize} \quad \frac{E_{x_1}^{\pi} \left[\sum_{n=1}^{\infty} \beta^{n-1} r_n \right]}{E_{x_1}^{\pi} \left[\sum_{n=1}^{\infty} \beta^{n-1} R_n \right]} \quad \text{subject to} \quad (i) \quad \pi \in \Pi \quad (8)$$

where

$$\begin{aligned}\sum_{n=1}^{\infty} \beta^{n-1} r_n &= r(X_1, U_1) + \beta r(X_2, U_2) + \cdots + \beta^{n-1} r(X_n, U_n) + \cdots \\ \sum_{n=1}^{\infty} \beta^{n-1} R_n &= R(X_1, U_1) + \beta R(X_2, U_2) + \cdots + \beta^{n-1} R(X_n, U_n) + \cdots\end{aligned}$$

Here Π is the set of all Markov policies, whose element $\pi = \{\pi_1, \pi_2, \dots, \pi_n, \dots\}$ is an infinite sequence of decision functions :

$$\pi_n : S \rightarrow A \quad n = 1, 2, \dots \quad (9)$$

An introduction of Lagrange multiplier λ reduces the fractional optimization problem (8) to a standard discounted dynamic programming problem ([3],[5],[6],[9]) as follows:

$$\begin{aligned} & \text{Maximize} \quad E_{x_1}^{\pi} \left[\sum_{n=1}^{\infty} \beta^{n-1} (r_n - \lambda R_n) \right] \quad (10) \\ P'(x_1; \lambda) \quad & \text{subject to} \quad \text{(i) } x_{n+1} \sim p(\cdot | x_n, u_n) \quad n = 1, 2, \dots \\ & \quad \quad \quad \text{(ii) } u_n \in A \quad n = 1, 2, \dots \\ & \quad \quad \quad x_1 \in S, \quad \lambda \in R^1.\end{aligned}$$

Let $u(x_1; \lambda)$ be the maximum value of the problem (10). Then we have the recursive equation:

THEOREM 3.1

$$\begin{aligned} u(x; \lambda) &= \text{Max}_{u \in A} [r(x, u) - \lambda R(x, u) + \beta \sum_{y \in S} u(y; \lambda) p(y | x, u)] \quad (11) \\ & \quad \quad \quad x \in S, \quad \lambda \in R^1.\end{aligned}$$

4 Fractional Programming Approach

In this section we solve the fractional expectation problems (4) and (8) through both fractional programming and dynamic programming.

4.1 Fractional Programming

Let us review two fundamental results on fractional programming. We consider the following problem:

$$\text{Fr} \quad \text{Maximize} \quad \frac{f(z)}{g(z)} \quad \text{subject to} \quad z \in Z \quad (12)$$

where Z is a nonempty set and $f : Z \rightarrow R^1$, $g : Z \rightarrow (0, \infty)$. It is well-known that the fractional programming problem Fr is associated with the following parametric problem:

$$\text{Pr}(\lambda) \quad \text{Maximize} \quad f(z) - \lambda g(z) \quad \text{subject to} \quad z \in Z. \quad (13)$$

THEOREM 4.1 ([11]) *The fractional problem Fr has an optimal solution $z^* \in Z$ if and only if the parametric problem $\text{Pr}(\lambda)$ has the optimal solution $z^* \in Z$ for some parameter λ and the optimal value vanishes.*

Let us consider Dinkelbach's Algorithm:

- Step 1. Select some $z \in Z$ and set $n = 1$, $z_{(1)} = z$ and $\lambda_{(1)} = \frac{f(z)}{g(z)}$.
- Step 2. Solve $\text{Pr}(\lambda_{(n)})$ and select some optimal solution $z \in Z$.
- Step 3. If $f(z) - \lambda_{(n)}g(z) = 0$, set $z' = z$ and $\lambda' = \frac{f(z)}{g(z)}$, and stop. Otherwise, set $z_{(n+1)} = z$ and $\lambda_{(n+1)} = \frac{f(z)}{g(z)}$.
- Step 4. Set $n = n + 1$ and go to Step 2.

THEOREM 4.2 ([11]) *Either Dinkelbach's Algorithm terminates in some finite n -th iteration, in which case z' is an optimal solution and λ' is a maximum value of Fr, or else the sequence $\{\lambda_{(n)}\}$ converges strict-monotonically to the maximum value of Fr. Termination is assured if Z is finite.*

We remark that the convergence is in fact superlinear. If Dinkelbach's Algorithm generates a finite sequence $\{\lambda_{(k)}\}_{1 \leq k \leq n}$ with properties

$$(i) \quad \lambda_{(1)} < \lambda_{(2)} < \cdots < \lambda_{(n-1)} < \lambda_{(n)},$$

- (ii) $f(z) - \lambda_{(n)}g(z) = 0$ for some optimal solution $z \in Z$ of $\text{Pr}(\lambda_{(n)})$, (iii) $z' = z$, and (iv) $\lambda' = \frac{f(z)}{g(z)}$, and terminates, then the z is an optimal solution and $\lambda_{(n)}$ is the maximum value of Fr.

4.2 Fractional Expectation Problems

First let us consider the fractional expectation problem (4) by use of fractional programming ([1]) and dynamic programming. The problem (4) is formulated as the following fractional programming problem:

$$\text{Fr}(x_1) \quad \text{Maximize} \quad \frac{f(\pi; x_1)}{g(\pi; x_1)} \quad \text{subject to} \quad \pi \in \Pi \quad (14)$$

where Π is the set of N -stage Markov policies and

$$\begin{aligned} f(\pi; x_1) &= E_{x_1}^{\pi} \left[\sum_{n=1}^{N+1} r_n \right] \\ g(\pi; x_1) &= E_{x_1}^{\pi} \left[\sum_{n=1}^{N+1} R_n \right]. \end{aligned}$$

Then the corresponding parametric problem reduces to:

$$\text{Pr}(x_1)(\lambda) \quad \text{Maximize} \quad f(\pi; x_1) - \lambda g(\pi; x_1) \quad \text{subject to} \quad \pi \in \Pi. \quad (15)$$

THEOREM 4.3 For each initial state $x_1 \in X$, Dinkelbach's Algorithm yields a Markov policy π^* , which is optimal at x_1 :

$$\frac{E_{x_1}^{\pi^*} \left[\sum_{n=1}^{N+1} r_n \right]}{E_{x_1}^{\pi^*} \left[\sum_{n=1}^{N+1} R_n \right]} \geq \frac{E_{x_1}^{\pi} \left[\sum_{n=1}^{N+1} r_n \right]}{E_{x_1}^{\pi} \left[\sum_{n=1}^{N+1} R_n \right]} \quad \forall \pi \in \Pi. \quad (16)$$

Proof Since Π is finite, Theorems 4.1 and 4.2 apply. \square

Second we consider the infinite-stage problem (8). By taking in turn

$$\begin{aligned} f(\pi; x_1) &= E_{x_1}^{\pi} \left[\sum_{n=1}^{\infty} \beta^{n-1} r_n \right] \\ g(\pi; x_1) &= E_{x_1}^{\pi} \left[\sum_{n=1}^{\infty} \beta^{n-1} R_n \right], \end{aligned}$$

we have a stationary policy which is optimal at a given initial state.

THEOREM 4.4 For each state $x_1 \in X$, Dinkelbach's Algorithm yields a stationary policy $\pi^* = h^{(\infty)}$, which is optimal at x_1 :

$$\frac{E_{x_1}^{\pi^*} \left[\sum_{n=1}^{\infty} \beta^{n-1} r_n \right]}{E_{x_1}^{\pi^*} \left[\sum_{n=1}^{\infty} \beta^{n-1} R_n \right]} \geq \frac{E_{x_1}^{\pi} \left[\sum_{n=1}^{\infty} \beta^{n-1} r_n \right]}{E_{x_1}^{\pi} \left[\sum_{n=1}^{\infty} \beta^{n-1} R_n \right]} \quad \forall \pi \in \Pi \quad (17)$$

where $h : S \rightarrow A$ is a stage-free decision function of π^* :

$$h^{(\infty)} = \{h, h, \dots, h, \dots\}.$$

Proof Let Π_{st} be the set of all stationary policies. Then we see that $\Pi_{st} \subset \Pi$ and Π_{st} is finite. We restrict the fractional problem (14) to Π_{st} . Then Theorems 4.1 and 4.2 apply. In fact, the corresponding parametric problem (15) is a discounted dynamic programming problem in the sense of D. Blackwell ([3]). Thus it has an optimal stationary policy. \square

5 A 2-2 Decision Models

In this section, we illustrate a two-state and two-action decision model.

5.1 A 2-2-2 Decision Model

As an illustrative example we consider the following two-stage problem:

$$\begin{aligned} & \text{Maximize} && \frac{E_{x_1}^{\pi} [r(x_1, u_1) + r(X_2, U_2) + k(X_3)]}{E_{x_1}^{\pi} [R(x_1, u_1) + R(X_2, U_2) + K(X_3)]} \\ F(x_1) & \text{subject to} && \text{(i) } x_{n+1} \sim p(\cdot | x_n, u_n) \quad 1 \leq n \leq 2 \\ & && \text{(ii) } u_n \in A \quad 1 \leq n \leq 2 \end{aligned} \quad (18)$$

on the following data:

stage rewards : $(r(x_t, u_t), R(x_t, u_t))$			terminal rewards	
$x_t \setminus u_t$	a_1	a_2	x_3	$(k(x_3), K(x_3))$
s_1	(0, 2)	(1, 1)	s_1	(1, 2)
s_2	(-1, 3)	(2, 2)	s_2	(0, 1)

transition law					
$P(a_1) = \{p(x_{t+1} x_t, a_1)\}$		$P(a_2) = \{p(x_{t+1} x_t, a_2)\}$			
$x_t \setminus x_{t+1}$	s_1	s_2	$x_t \setminus x_{t+1}$	s_1	s_2
s_1	1/2	1/2	s_1	1	0
s_2	0	1	s_2	1/4	3/4

Thus we have the following parametric data:

stage reward : $r(x_t, u_t) - \lambda R(x_t, u_t)$			terminal reward	
$x_t \setminus u_t$	a_1	a_2	x_3	$k(x_3) - \lambda K(x_3)$
s_1	$0 - 2\lambda$	$1 - \lambda$	s_1	$1 - 2\lambda$
s_2	$-1 - 3\lambda$	$2 - 2\lambda$	s_2	$0 - \lambda$

Then the recursive equation

$$\begin{aligned}
 u_3(x; \lambda) &= k(x) - \lambda K(x) \\
 u_2(x; \lambda) &= \text{Max}_{u \in A} \left[r(x, u) - \lambda R(x, u) + \sum_{y \in S} u_3(y; \lambda) p(y|x, u) \right] \\
 u_1(x; \lambda) &= \text{Max}_{u \in A} \left[r(x, u) - \lambda R(x, u) + \sum_{y \in S} u_2(y; \lambda) p(y|x, u) \right] \\
 & \quad x \in S, \lambda \in R^1
 \end{aligned} \tag{19}$$

together with the suffixed notations

$$u_n(\lambda) := u_n(s_1; \lambda), \quad v_n(\lambda) := u_n(s_2; \lambda)$$

$$k_i := k(s_i), \quad K_i := K(s_i), \quad r_i^k := r(s_i, a_k), \quad R_i^k := R(s_i, a_k), \quad p_{ij}^k := p(s_j | s_i, a_k)$$

reduces to:

$$\begin{aligned}
 u_3(\lambda) &= k_1 - \lambda K_1 \\
 v_3(\lambda) &= k_2 - \lambda K_2 \\
 u_n(\lambda) &= \left[r_1^1 - \lambda R_1^1 + p_{11}^1 u_{n+1}(\lambda) + p_{12}^1 v_{n+1}(\lambda) \right] \\
 & \quad \vee \left[r_1^2 - \lambda R_1^2 + p_{11}^2 u_{n+1}(\lambda) + p_{12}^2 v_{n+1}(\lambda) \right] \\
 v_n(\lambda) &= \left[r_2^1 - \lambda R_2^1 + p_{21}^1 u_{n+1}(\lambda) + p_{22}^1 v_{n+1}(\lambda) \right] \\
 & \quad \vee \left[r_2^2 - \lambda R_2^2 + p_{21}^2 u_{n+1}(\lambda) + p_{22}^2 v_{n+1}(\lambda) \right] \quad n = 1, 2.
 \end{aligned} \tag{20}$$

Then Eq.(20) becomes:

$$\begin{aligned} u_3(\lambda) &= 1 - 2\lambda \\ v_3(\lambda) &= 0 - \lambda \\ u_n(\lambda) &= \left[0 - 2\lambda + \frac{1}{2}u_{n+1}(\lambda) + \frac{1}{2}v_{n+1}(\lambda)\right] \vee [1 - \lambda + u_{n+1}(\lambda)] \\ v_n(\lambda) &= [-1 - 3\lambda + v_{n+1}(\lambda)] \vee \left[2 - 2\lambda + \frac{1}{4}u_{n+1}(\lambda) + \frac{3}{4}v_{n+1}(\lambda)\right] \quad n = 1, 2. \end{aligned}$$

Thus we have

$$\begin{aligned} u_2(\lambda) &= \left[\frac{1}{2} - \frac{7}{2}\lambda\right] \vee [2 - 3\lambda] = \begin{cases} \frac{1}{2} - \frac{7}{2}\lambda, & -\infty < \lambda \leq -3 \\ 2 - 3\lambda, & -3 \leq \lambda < \infty \end{cases} \\ v_2(\lambda) &= [-1 - 4\lambda] \vee \left[\frac{9}{4} - \frac{13}{4}\lambda\right] = \begin{cases} -1 - 4\lambda, & -\infty < \lambda \leq -\frac{13}{3} \\ \frac{9}{4} - \frac{13}{4}\lambda, & -\frac{13}{3} \leq \lambda < \infty \end{cases} \\ u_1(\lambda) &= \begin{cases} -\frac{1}{4} - \frac{23}{4}\lambda, & -\infty < \lambda \leq -\frac{13}{3} \\ \frac{11}{8} - \frac{43}{8}\lambda, & -\frac{13}{3} \leq \lambda \leq -3 \\ \frac{17}{8} - \frac{41}{8}\lambda, & -3 \leq \lambda \leq -\frac{7}{9} \\ 3 - 4\lambda, & -\frac{7}{9} \leq \lambda < \infty \end{cases} \\ v_1(\lambda) &= \begin{cases} -2 - 7\lambda, & -\infty < \lambda \leq -\frac{13}{3} \\ \frac{5}{4} - \frac{25}{4}\lambda, & -\frac{13}{3} \leq \lambda \leq -\frac{47}{17} \\ \frac{67}{16} - \frac{83}{16}\lambda, & -\frac{47}{17} \leq \lambda < \infty \end{cases} \end{aligned}$$

Then the desired optimal policy $\pi^*(\lambda) = \{\pi_1^*(\lambda), \pi_2^*(\lambda)\}$ where

$$\pi_n^*(\lambda) = \begin{bmatrix} \pi_n^*(s_1; \lambda) \\ \pi_n^*(s_2; \lambda) \end{bmatrix} \quad (21)$$

is specified as follows :

$$\pi_2^*(\lambda) = \begin{cases} \begin{bmatrix} a_1 \\ a_1 \end{bmatrix}, & -\infty < \lambda \leq -\frac{13}{3} \\ \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}, & -\frac{13}{3} \leq \lambda \leq -3 \\ \begin{bmatrix} a_2 \\ a_2 \end{bmatrix}, & -3 \leq \lambda < \infty \end{cases} \quad \pi_1^*(\lambda) = \begin{cases} \begin{bmatrix} a_1 \\ a_1 \end{bmatrix}, & -\infty < \lambda \leq -\frac{47}{17} \\ \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}, & -\frac{47}{17} \leq \lambda \leq -\frac{7}{9} \\ \begin{bmatrix} a_2 \\ a_2 \end{bmatrix}, & -\frac{7}{9} \leq \lambda < \infty \end{cases} \quad (22)$$

By applications of Dinkelbach's Algorithm from $\pi = \left\{ \begin{bmatrix} a_1 \\ a_1 \end{bmatrix}, \begin{bmatrix} a_1 \\ a_1 \end{bmatrix} \right\}$, we have optimal solutions as follows:

CASE(I) Algorithm I for $x_1 = s_1$.

1. Select $\pi_1 = \left\{ \left[\begin{array}{c} a_1 \\ a_1 \end{array} \right], \left[\begin{array}{c} a_1 \\ a_1 \end{array} \right] \right\} \in \Pi$. Then $\lambda_{(1)} = \frac{f(\pi_1; s_1)}{g(\pi_1; s_1)} = \frac{-1/4}{23/4} = -\frac{1}{23}$.
2. Solve $\Pr\left(-\frac{1}{23}\right)$ and select unique optimal solution $\pi_2 = \left\{ \left[\begin{array}{c} a_2 \\ a_2 \end{array} \right], \left[\begin{array}{c} a_2 \\ a_2 \end{array} \right] \right\} \in \Pi$. Then $f(\pi_2; s_1) - \lambda_{(1)}g(\pi_2; s_1) = 3 - \left(-\frac{1}{23}\right) \cdot 4 = \frac{72}{23} \neq 0$. Hence $\lambda_{(2)} = \frac{f(\pi_2; s_1)}{g(\pi_2; s_1)} = \frac{3}{4}$.
3. Solve $\Pr\left(\frac{3}{4}\right)$ and select unique optimal solution $\pi^* = \left\{ \left[\begin{array}{c} a_2 \\ a_2 \end{array} \right], \left[\begin{array}{c} a_2 \\ a_2 \end{array} \right] \right\} \in \Pi$. Then $f(\pi^*; s_1) - \lambda_{(2)}g(\pi^*; s_1) = 3 - \frac{3}{4} \cdot 4 = 0$. Thus $\pi^* = \pi_2$ is an optimal at s_1 and $\lambda_{(2)} = \frac{3}{4}$ is the desired maximum value.

CASE(II) Algorithm I for $x_1 = s_2$.

1. Select $\pi_1 = \left\{ \left[\begin{array}{c} a_1 \\ a_1 \end{array} \right], \left[\begin{array}{c} a_1 \\ a_1 \end{array} \right] \right\} \in \Pi$. Then $\lambda_{(1)} = \frac{f(\pi_1; s_2)}{g(\pi_1; s_2)} = \frac{-2}{7}$.
2. Solve $\Pr\left(-\frac{2}{7}\right)$ and select unique optimal solution $\pi_2 = \left\{ \left[\begin{array}{c} a_2 \\ a_2 \end{array} \right], \left[\begin{array}{c} a_2 \\ a_2 \end{array} \right] \right\} \in \Pi$. Then $f(\pi_2; s_2) - \lambda_{(1)}g(\pi_2; s_2) = \frac{67}{16} - \left(-\frac{2}{7}\right) \cdot \frac{83}{16} = \frac{635}{112} \neq 0$. Hence $\lambda_{(2)} = \frac{f(\pi_2; s_2)}{g(\pi_2; s_2)} = \frac{67/16}{83/16} = \frac{67}{83}$.
3. Solve $\Pr\left(\frac{67}{83}\right)$ and select unique optimal solution $\pi^* = \left\{ \left[\begin{array}{c} a_2 \\ a_2 \end{array} \right], \left[\begin{array}{c} a_2 \\ a_2 \end{array} \right] \right\} \in \Pi$. Then $f(\pi^*; s_2) - \lambda_{(2)}g(\pi^*; s_2) = \frac{67}{16} - \frac{67}{83} \cdot \frac{83}{16} = 0$. Thus $\pi^* = \pi_2$ is also optimal at s_2 and $\lambda_{(2)} = \frac{67}{83}$ is the desired maximum value.

Therefore, the resulting stationary policy $\pi^* = \left\{ \left[\begin{array}{c} a_2 \\ a_2 \end{array} \right], \left[\begin{array}{c} a_2 \\ a_2 \end{array} \right] \right\}$ is optimal (for both states) and the optimal ratio vectors is $\left(\begin{array}{c} 3/4 \\ 67/83 \end{array} \right)$.

5.2 A 2-2- ∞ Decision Model

Now we consider the corresponding infinite-stage problem on the two-state and two-action model:

$$F'(x_1) \quad \text{Maximize} \quad \frac{E_{x_1}^{\pi} \left[\sum_{n=1}^{\infty} \beta^{n-1} r_n \right]}{E_{x_1}^{\pi} \left[\sum_{n=1}^{\infty} \beta^{n-1} R_n \right]} \quad \text{subject to} \quad (i) \quad \pi \in \Pi \quad (23)$$

where $\beta = 0.8$. Then the recursive equation for the corresponding parametric problem

$$u(x; \lambda) = \text{Max}_{u \in A} \left[r(x, u) - \lambda R(x, u) + \beta \sum_{y \in S} u(y; \lambda) p(y|x, u) \right] \quad (24)$$

$x \in S, \lambda \in R^1$

together with the suffixed notations

$$u(\lambda) := u(s_1; \lambda), \quad v(\lambda) := u(s_2; \lambda)$$

$$r_i^k := r(s_i, a_k), \quad R_i^k := R(s_i, a_k), \quad p_{ij}^k := p(s_j | s_i, a_k)$$

reduces to:

$$u(\lambda) = \left[r_1^1 - \lambda R_1^1 + \beta(p_{11}^1 u(\lambda) + p_{12}^1 v(\lambda)) \right] \vee \left[r_1^2 - \lambda R_1^2 + \beta(p_{11}^2 u(\lambda) + p_{12}^2 v(\lambda)) \right] \quad (25)$$

$$v(\lambda) = \left[r_2^1 - \lambda R_2^1 + \beta(p_{21}^1 u(\lambda) + p_{22}^1 v(\lambda)) \right] \vee \left[r_2^2 - \lambda R_2^2 + \beta(p_{21}^2 u(\lambda) + p_{22}^2 v(\lambda)) \right].$$

Then Eq.(25) reduces to:

$$u(\lambda) = \left[0 - 2\lambda + \frac{4}{5} \left(\frac{1}{2}u(\lambda) + \frac{1}{2}v(\lambda) \right) \right] \vee \left[1 - \lambda + \frac{4}{5}u(\lambda) \right]$$

$$v(\lambda) = \left[-1 - 3\lambda + \frac{4}{5}v(\lambda) \right] \vee \left[2 - 2\lambda + \frac{4}{5} \left(\frac{1}{4}u(\lambda) + \frac{3}{4}v(\lambda) \right) \right]$$

namely

$$[-10\lambda - 3u(\lambda) + 2v(\lambda)] \vee [5 - 5\lambda - u(\lambda)] = 0$$

$$[-5 - 15\lambda - v(\lambda)] \vee [10 - 10\lambda + u(\lambda) - 2v(\lambda)] = 0.$$

This system of two function equations has the following unique solution:

$$u(\lambda) = \begin{cases} -\frac{10}{3} - \frac{40}{3}\lambda, & -\infty < \lambda \leq -\frac{5}{2} \\ 5 - 10\lambda, & -\frac{5}{2} \leq \lambda \leq 0 \\ 5 - 5\lambda, & 0 \leq \lambda < \infty \end{cases} \quad v(\lambda) = \begin{cases} -5 - 15\lambda, & -\infty < \lambda \leq -\frac{5}{2} \\ \frac{15}{2} - 10\lambda, & -\frac{5}{2} \leq \lambda \leq 0 \\ \frac{15}{2} - \frac{15}{2}\lambda, & 0 \leq \lambda < \infty. \end{cases}$$

Then the desired optimal policy $\pi^*(\lambda) = h^{(\infty)}(\lambda)$ where

$$h(\lambda) = \begin{bmatrix} h(s_1; \lambda) \\ h(s_2; \lambda) \end{bmatrix} \quad (26)$$

is specified as follows :

$$h(\lambda) = \begin{cases} \begin{bmatrix} a_1 \\ a_1 \end{bmatrix}, & -\infty < \lambda \leq -\frac{5}{2} \\ \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}, & -\frac{5}{2} \leq \lambda \leq 0 \\ \begin{bmatrix} a_2 \\ a_2 \end{bmatrix}, & 0 \leq \lambda < \infty \end{cases} \quad (27)$$

By applications of Dinkelbach's Algorithm from $\pi = h^{(\infty)}$ with $h = \begin{bmatrix} a_1 \\ a_1 \end{bmatrix}$, we have the following optimal solutions:

CASE(I) Algorithm II for $x_1 = s_1$.

1. Select $\pi_1 = h_1^{(\infty)} \in \Pi_{st}$ with $h_1 = \begin{bmatrix} a_1 \\ a_1 \end{bmatrix}$. Then $\lambda_{(1)} = \frac{f(\pi_1; s_1)}{g(\pi_1; s_1)} = \frac{-15/3}{-40/3} = -\frac{3}{8}$.
2. Solve $\Pr\left(-\frac{3}{8}\right)$ and select optimal solution $\pi_2 = h_2^{(\infty)} \in \Pi_{st}$ with $h_2 = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$. Then $f(\pi_2; s_1) - \lambda_{(1)}g(\pi_2; s_1) = 5 - \left(-\frac{3}{8}\right) \cdot 10 = \frac{35}{4} \neq 0$. Hence $\lambda_{(2)} = \frac{f(\pi_2; s_1)}{g(\pi_2; s_1)} = \frac{5}{10} = \frac{1}{2}$.
3. Solve $\Pr\left(\frac{1}{2}\right)$ and select optimal solution $\pi_3 = h_3^{(\infty)} \in \Pi_{st}$ with $h_3 = \begin{bmatrix} a_2 \\ a_2 \end{bmatrix}$. Then $f(\pi_3; s_1) - \lambda_{(2)}g(\pi_3; s_1) = 5 - \frac{1}{2} \cdot 5 = \frac{5}{2} \neq 0$. Hence $\lambda_{(3)} = \frac{f(\pi_3; s_1)}{g(\pi_3; s_1)} = \frac{5}{5} = 1$.
4. Solve $\Pr(1)$ and select optimal solution $\pi^* = h_*^{(\infty)} \in \Pi_{st}$ with $h_* = \begin{bmatrix} a_2 \\ a_2 \end{bmatrix}$.
5. Then $f(\pi^*; s_1) - \lambda_{(3)}g(\pi^*; s_1) = 5 - 1 \cdot 5 = 0$. Thus $\pi^* = \pi_3$ is an optimal at s_1 and $\lambda_{(3)} = 1$ is the desired maximum value.

CASE(II) Algorithm II for $x_1 = s_2$.

1. Select $\pi_1 = h_1^{(\infty)} \in \Pi_{st}$ with $h_1 = \begin{bmatrix} a_1 \\ a_1 \end{bmatrix}$. Then $\lambda_{(1)} = \frac{f(\pi_1; s_2)}{g(\pi_1; s_2)} = \frac{-5}{15} = -\frac{1}{3}$.
2. Solve $\Pr\left(-\frac{1}{3}\right)$ and select optimal solution $\pi_2 = h_2^{(\infty)} \in \Pi_{st}$ with $h_2 = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$. Then $f(\pi_2; s_2) - \lambda_{(1)}g(\pi_2; s_2) = \frac{15}{2} - \left(-\frac{1}{3}\right) \cdot 10 = \frac{65}{6} \neq 0$. Hence $\lambda_{(2)} = \frac{f(\pi_2; s_2)}{g(\pi_2; s_2)} = \frac{15/2}{10} = \frac{3}{4}$.
3. Solve $\Pr\left(\frac{1}{2}\right)$ and select optimal solution $\pi_3 = h_3^{(\infty)} \in \Pi_{st}$ with $h_3 = \begin{bmatrix} a_2 \\ a_2 \end{bmatrix}$. Then $f(\pi_3; s_2) - \lambda_{(2)}g(\pi_3; s_2) = \frac{15}{2} - \frac{3}{4} \cdot \frac{15}{2} = \frac{15}{8} \neq 0$. Hence $\lambda_{(3)} = \frac{f(\pi_3; s_2)}{g(\pi_3; s_2)} = \frac{15/2}{15/2} = 1$.
4. Solve $\Pr(1)$ and select optimal solution $\pi^* = h_*^{(\infty)} \in \Pi_{st}$ with $h_* = \begin{bmatrix} a_2 \\ a_2 \end{bmatrix}$.
5. Then $f(\pi^*; s_2) - \lambda_{(3)}g(\pi^*; s_2) = \frac{15}{2} - 1 \cdot \frac{15}{2} = 0$. Thus $\pi^* = \pi_3$ is also an optimal at s_2 and $\lambda_{(3)} = 1$ is also the desired maximum value.

On the other hand, applications from $\pi = h^{(\infty)}$ with $h = \begin{bmatrix} a_2 \\ a_1 \end{bmatrix}$ yields the following results:

CASE(III) Algorithm II for $x_1 = s_1$.

1. Select $\pi_1 = h_1^{(\infty)} \in \Pi_{st}$ with $h_1 = \begin{bmatrix} a_2 \\ a_1 \end{bmatrix}$. Then $\lambda_{(1)} = \frac{f(\pi_1; s_1)}{g(\pi_1; s_1)} = \frac{5}{5} = 1$. From CASE(I), the desired maximum value is 1. Thus the policy π_1 is also optimal at s_1 .
2. Solve $\Pr(1)$ and select optimal solution $\pi_2 = h_2^{(\infty)} \in \Pi_{st}$ with $h_2 = \begin{bmatrix} a_2 \\ a_2 \end{bmatrix}$. Then $f(\pi_2; s_1) - \lambda_{(1)}g(\pi_2; s_1) = 5 - 1 \cdot 5 = 0$. Thus $\pi^* = \pi_2$ is optimal at s_1 and $\lambda_{(2)} = 1$ is the desired maximum value.

CASE(IV) Algorithm II for $x_1 = s_2$.

1. Select $\pi_1 = h_1^{(\infty)} \in \Pi_{st}$ with $h_1 = \begin{bmatrix} a_2 \\ a_1 \end{bmatrix}$. Then $\lambda_{(1)} = \frac{f(\pi_1; s_2)}{g(\pi_1; s_2)} = \frac{-5}{15} = -\frac{1}{3}$.

2. Solve $\Pr\left(-\frac{1}{3}\right)$ and select unique optimal solution $\pi_2 = h_2^{(\infty)} \in \Pi_{st}$ with $h_2 = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$.

Hereafter CASE(II) follows. Thus, $\pi^* = \pi_3$ is also an optimal at s_2 and $\lambda_{(3)} = 1$ is also the desired maximum value. Thus the policy π_1 is not optimal at s_2 .

Therefore, the resulting stationary policy $\pi^* = h_1^{(\infty)}$ with $h_1 = \begin{bmatrix} a_2 \\ a_2 \end{bmatrix}$ is optimal (for both states) and the optimal ratio vectors is $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$. Furthermore, the stationary policy $\pi^{**} = h_2^{(\infty)}$ with $h_2 = \begin{bmatrix} a_2 \\ a_1 \end{bmatrix}$ is optimal at s_2 .

References

- [1] A.I. Barros, *Discrete and Fractional Programming Techniques for Location Models*, Kluwer Academic Publishers, Dordrecht, 1998.
- [2] R.E. Bellman, *Dynamic Programming*, Princeton Univ. Press, NJ, 1957.
- [3] D. Blackwell, Discounted dynamic programming, *Ann. Math. Stat.* **36**(1965), 226-235.
- [4] T. Fujita, Re-examination of Markov Policies for Additive Decision Process, *Bull. Infor. Cyber.*, 29(1997), 51-65.
- [5] K. Hinderer, *Foundations of Non-Stationary Dynamic Programming with Discrete Time Parameter*, Lect. Notes in Operation Research and Mathematical Systems, Vol. 33, Springer-Verlag, Berlin, 1970.
- [6] R. A. Howard, *Dynamic Programming and Markov Processes*, MIT Press, Cambridge, Mass., 1960.
- [7] S. Iwamoto, On expected values of Markov statistics, *Bull. Infor. Cyber.*, 30(1998), 1-24.
- [8] S. Iwamoto and T. Fujita, On conditional expected values of Markov statistics, under preparation.
- [9] M.L. Puterman, *Markov Decision Processes : discrete stochastic dynamic programming*, Wiley & Sons, New York, 1994.
- [10] M. Sniedovich, Analysis of a class of fractional programming problems, *Math. Prog.* **43**(1989), 329-347.
- [11] M. Sniedovich, *Dynamic Programming*, Marcel Dekker, Inc. NY, 1992.