

Stochastic shortest path problems with associative accumulative criteria

Yoshio Ohtsubo

*Department of Mathematics, Faculty of Science,
Kochi University, Kochi 780-8520, Japan*

Abstract We consider a stochastic shortest path problem with associative criteria in which for each node of a graph we choose a probability distribution over the set of successor nodes so as to reach a given target node optimally. We formulate such a problem as an associative Markov decision processes. We show that an optimal value function is a unique solution to an optimality equation and find an optimal stationary policy. Also we give a value iteration method and a policy improvement method.

Keywords: shortest path problem; Markov decision process; associative criterion; invariant imbedding method; optimality equation; existence of optimal policy

1. Introduction.

For a directed graph with nodes $1, 2, \dots, K$ and with a cost (length or time) assigned to each arc, a stochastic shortest path problem is to select a probability distribution over all possible successor nodes at each node $i \neq K$ so as to reach a target node K with minimal associative accumulative cost.

Such a stochastic shortest path problem is analyzed by using the general theory of Markov decision processes in many references. Eaton and Zadeh[3] formulated such a problem as a pursuit problem and they showed that the optimal expected total cost is a unique solution to an optimality equation if at least one proper policy exists, and they gave an optimal value by a value iteration method. Derman in [4, 5] considered the problem, where a target state (node) is absorbing, and proved that the problem has an optimal stationary policy and he gave several methods for obtaining optimal solutions. In [16], Sancho formulated Markov decision processes to analyze the problem and gave a policy iteration method. Bertsekas and Tsitsiklis[2] investigated the problem without the cost nonnegativity assumption and proved a natural generalization of the standard result for the deterministic shortest path problem within the framework of undiscounted finite state Markovian decision processes. In all of these, a criterion function is the expected total cost, which we call an additive case.

Also, Ohtsubo[12] considered a minimizing risk models in stochastic shortest path problems as undiscounted finite Markov processes and showed that an optimal value function is a unique solution to an optimality equation and found an optimal stationary policy by using an invariant imbedding method. General minimizing risk models in discounted Markov decision processes were investigated in White[18], Wu and Lin[19], Ohtsubo and Toyonaga[11, 13] and Ohtsubo[14].

On the other hand, Maruyama in [9, 10] investigated deterministic shortest path problems with associative criteria and show the existence and uniqueness of the optimal value. Especially in [10] he obtained a parameterized recursive equation for the class of the problem by using an invariant imbedding technique.

Furthermore the optimization problem for minimum criteria, which is associative, was first introduced by Bellman and Zadeh[1] as decision-making in fuzzy environment, and Iwamoto et al.[6, 7, 8] and Ohtsubo[15] formulated their optimization problem as finite horizon Markov decision processes and gave a optimal policy by using an invariant imbedding approach.

In this paper we concern ourselves with a stochastic shortest path problem with an associative criterion, which is an expected accumulate cost $E_i^\pi[\bigcirc_{n=1}^\tau Y_n] = E_i^\pi[Y_0 \circ Y_1 \circ Y_2 \circ \dots \circ Y_\tau]$ where Y_n is a cost at n th step, \circ is an operator with an associative property satisfying some conditions, τ is a hitting time to the target node K and E_i^π is an expectation operator when the starting node is i and a policy π is used. In Section 2, we give notations and formulate our model as undiscounted finite Markov decision processes with infinite horizon. In Section 3, we prove that the optimal value function is a unique solution to an optimality equation by using an invariant imbedding approach and that it is given by a value iteration method. We also show that there exists an optimal left continuous stationary policy. In Section 4, we

E-mail address : ohtsubo@math.kochi-u.ac.jp

give a policy improvement method for obtained a optimal policy.

2. Notations and formulation

In this section we formulate associative models in stochastic shortest path problems as Markov decision Processes $\Gamma = ((X_n), (A_n), (Y_n), p)$ with a discrete time space $N = \{0, 1, 2, \dots\}$. The state space S is a finite set $\{1, 2, \dots, K\}$ where K is a target state, and we denote the state at time $n \in N$ by X_n . The action space A is finite and we denote the action at time $n \in N$ by A_n . The cost space E is a finite set $\{y_1, y_2, \dots, y_\ell\}$, where $E \subset B$ for some subset B of R , and $Y_n \in E$ is a random cost function at time $n \in N$ with $Y_0 = e$, where e is a unit element defined below. We define conditional probability distributions by

$$\begin{aligned} q^a(j|i) &= P(X_{n+1} = j | X_n = i, A_n = a), \\ \hat{q}_{ij}^a(y) &= P(Y_{n+1} = y | X_n = i, X_{n+1} = j, A_n = a) \end{aligned}$$

and set

$$p^a(j, y|i) = q^a(j|i)\hat{q}_{ij}^a(y) = P(X_{n+1} = j, Y_{n+1} = y | X_n = i, A_n = a)$$

for $i, j \in S, a \in A$ and $y \in E$. We use $S_B = S \times B$ as a new state space.

For a binary operator $\circ : R \times R \rightarrow R$ and a subset B of R , we assume that

- (i) B is closed for the operator \circ , that is, $x \circ y \in B$ for any $x, y \in B$,
- (ii) the operator \circ is associative, that is, $(x \circ y) \circ z = x \circ (y \circ z)$ ($= x \circ y \circ z$, say) for any $x, y, z \in B$,
- (iii) B has a unit element e , that is, $e \in B$ and $x \circ e = e \circ x = x$ for any $x \in B$,
- (iv) (B, \circ) is nondecreasing in the sense that $x \leq x \circ y$ and $x \leq y \circ x$ for any $x, y \in B$.

On the condition (iv), letting $x = e$, we notice that $y \geq e$ for any $y \in B$. Also we easily see under the conditions (i), (ii) and (iii) that if $x \geq e$ for any $x \in B$ and if $x \circ y \leq x \circ z$ and $y \circ x \leq z \circ x$ for any $x, y, z \in B$ satisfying $y \leq z$, then the condition (iv) holds. In algebra, (B, \circ) satisfying the conditions (i), (ii) and (iii) is called a semigroup and it is also analogous to t -conorm (or s -norm) in fuzzy set theory (cf. Zimmermann[20]).

We give several examples in which (B, \circ) satisfies the above conditions (cf. Maruyama[9] and [20]).

Example 2.1.

(1) (Additive case). When $x \circ y = \min(x+y-L, M)$ for constants L, M such that $-\infty < L < M \leq \infty$, we have $B = [L, M]$ and $e = L$, where $B = [L, \infty)$ if $M = \infty$. If $L = 0$ and $M = \infty$, it is a usual additive case, and if $L = 0$ and $M = 1$, it is called a bounded sum in the fuzzy set theory.

(2) (Multiplicative case). When $x \circ y = Lxy$ for a constant $L > 0$, we have $B = [1/L, \infty)$ and $e = 1/L$.

(3) (Maximum case) When $x \circ y = \max(x, y)$ for $x, y \in [L, M]$ where constants L, M satisfy $L < M$, we have $B = [L, M]$ and $e = L$.

(4) (Multiplicative-additive case). When $x \circ y = x+y-Lxy$ for a constant $L > 0$, we have $B = [0, 1/L]$ and $e = 0$.

(5) (Fractional case). When $x \circ y = (x+y)/(1+Lxy)$ for a constant $L > 0$, we have $B = [0, 1/\sqrt{L}]$ and $e = 0$. If $L = 1$, it is an Einstein sum in the fuzzy set theory.

(6) (Drastic maximum case) When $x \circ y = \begin{cases} \max(x, y) & \text{if } \min(x, y) = L \\ M & \text{otherwise} \end{cases}$ for $x, y \in [L, M]$ where L and M are constants so that $L < M$, we have $B = [L, M]$ and $e = L$.

(7) (Hamacher case) When $x \circ y = \frac{x+y-2xy}{1-xy}$ for each $x, y \in [0, 1)$, we have $B = [0, 1)$ and $e = 0$. This is a Hamacher sum in the fuzzy set theory.

Let a stopping time τ be a hitting time to the target state K , that is, τ is the smallest nonnegative integer n such that $X_n = K$, where $\tau = \infty$ if there does not exist such an integer n . Then we define the random reward as a criterion function by

$$Z = \bigcirc_{n=0}^{\tau} Y_n \equiv Y_0 \circ Y_1 \circ \dots \circ Y_\tau.$$

Then our problem is to minimize the expected reward $E_i^\pi[Z]$ with respect to all policies π .

To simplify the optimization problem, we can redefine the equivalent version of the Markov decision processes as follows. We assume that the target state K is absorbing and cost-free, that is, $q^a(K|K) = 1$ and $\hat{q}_{KK}^a(e) = 1$ and hence $p^a(K, e|K) = 1$ for all $a \in A$. Under this assumption we have

$$Z = \bigcirc_{k=0}^{\infty} Y_k \equiv \lim_{n \rightarrow \infty} \bigcirc_{k=0}^n Y_k,$$

which exists from the monotonicity of the assumption (iv), where we admit $Z = \infty$.

In order to analysis our problem we also define the random reward for a subproblem by

$$Z_n = \bigcirc_{k=0}^n Y_k \equiv Y_0 \circ Y_1 \circ \cdots \circ Y_n, \quad n \geq 0,$$

Further we define another random sequence as an imbedded parameter by

$$\Lambda_0 = \lambda, \quad \Lambda_{n+1} = \Lambda_n \circ Y_{n+1}, \quad n \geq 0,$$

where λ is a given initial parameter in B .

Let $H_0 = S_B$ and $H_{n+1} = H_n \times A \times S_B$ for each $n \in N$. Then H_n represents the set of all possible histories of the system when the n th action must be chosen, and we denote by θ_n the history at time $n \in N$. A decision rule δ_n for time $n \in N$ is a conditional probability given θ_n : $\delta_n(a_n|h_n) = P(A_n = a_n|\theta_n = h_n)$, where $h_n = (i_0, \lambda_0, a_0, i_1, \lambda_1, \dots, a_{n-1}, i_n, \lambda_n) \in H_n$ which is a realising value of $\theta_n = (X_0, \Lambda_0, A_0, X_1, \Lambda_1, \dots, A_{n-1}, X_n, \Lambda_n)$. It is assumed that $\delta_n(A_n \in A|h_n) = 1$ for every history $h_n = (i_0, \lambda_0, a_0, \dots, i_n, \lambda_n) \in H_n$. We denote by Δ the set of all decision rules. A policy π is an infinite sequence of decision rules $(\delta_n, n \geq 0) = (\delta_0, \delta_1, \delta_2, \dots, \delta_n, \dots)$. We denote by C the set of all such policies.

A policy $\pi = (\delta_n, n \geq 0)$ is said to be Markov when the decision rule δ_n is a function of $(X_n, \Lambda_n) = (i_n, \lambda_n)$ for every $n \in N$. We denote the set of such decision rules by Δ_M and the set of all Markov policies by C_M . Also, a policy π is called a deterministic Markov policy if π is Markov and $\delta_n(a|i, \lambda) = 1$ for some $a \in A$. We write $\delta_n(i, \lambda) = a$ for such a decision rule δ_n and we denote by Δ_D the set of such decision rules. We also denote the set of all deterministic Markov policies by C_D . When $\delta_n = \delta \in \Delta_D$ for all $n \in N$, we write $\pi = \delta^\infty$, which is called a stationary policy, and we denote the set of all stationary policies by C_D^s .

We denote by $E_i^\pi[Z]$ the conditional expectation of Z given an initial state $X_0 = i$ and a policy $\pi \in C$. Since the random variable Z depends upon not only i and π but also λ , we may sometimes use a conditional probability $P_{(i,\lambda)}^\pi(\cdot)$ and an expectation $E_{(i,\lambda)}^\pi(\cdot)$. Through this paper we assume that $P_{(i,\lambda)}^\pi(X_n = K \text{ for some } n \geq 0) = P_{(i,\lambda)}^\pi(\tau < \infty) = 1$ for every stationary policy $\pi \in C_D^s$ and each $(i, \lambda) \in S_B$, that is, the states $1, 2, \dots, K-1$ are transient when we use any policy $\pi \in C_D^s$. Thus we easily see that $P_{(i,\lambda)}^\pi(Z < \infty) = 1$ for all $\pi \in C_D^s$ and each $(i, \lambda) \in S_B$. This is analogous to a condition given in Ohtsubo[16].

A decision rule $\delta \in \Delta_D$ is said to be left continuous (on B) if for each $(i, \lambda) \in S_B$ there is a positive real number μ such that $\delta(i, \lambda) = \delta(i, \lambda - u)$ for all u such that $0 \leq u < \mu$ and $\lambda - u \in B$. A policy $\pi = \delta^\infty \in C_D^s$ is said to be left continuous if the decision rule δ is left continuous.

In order to analysis our model, we denote criterion functions for finite and infinite horizon cases by

$$F_n^\pi(i, \lambda) = E_i^\pi[\lambda \circ Z_n], \quad F^\pi(i, \lambda) = E_i^\pi[\lambda \circ Z],$$

respectively, for each $(i, \lambda) \in S_B$ and $\pi \in C$. When $n = 3$, the explicit form of the expectation $F_3^\pi(i_1, \lambda)$ is

$$\begin{aligned} E_{i_1}^\pi[\lambda \circ Z_3] &= \sum_{a_1 \in A} \sum_{y_1 \in E} \sum_{i_2 \in S} \sum_{a_2 \in A} \sum_{y_2 \in E} \sum_{i_3 \in S} \sum_{a_3 \in A} \sum_{y_3 \in E} \sum_{i_4 \in S} (\lambda \circ y_1 \circ y_2 \circ y_3) \\ &\quad \times p^{a_3}(i_4, y_3|i_3) \delta_2(a_3|i_1, \lambda, a_1, i_2, \lambda \circ y_1, a_2, i_3, \lambda \circ y_1 \circ y_2) \\ &\quad \times p^{a_2}(i_3, y_2|i_2) \delta_1(a_2|i_1, \lambda, a_1, i_2, \lambda \circ y_1) \\ &\quad \times p^{a_1}(i_2, y_1|i_1) \delta_0(a_1|i_1, \lambda) \end{aligned}$$

for $(i_1, \lambda) \in S_B$ and $\pi = (\delta_0, \delta_1, \delta_2, \dots) \in C$. We also define optimal value functions F_n^* and F^* for finite and infinite horizon cases by, respectively,

$$F_n^*(i, \lambda) = \inf_{\pi \in C} F_n^\pi(i, \lambda), \quad F^*(i, \lambda) = \inf_{\pi \in C} F^\pi(i, \lambda).$$

Then we notice that optimal value in the original problem is

$$F^*(i, e) = \sup_{\pi \in C} F^\pi(i, e) = \sup_{\pi \in C} E_i^\pi[Z],$$

since e is the unit element. A policy π is said to be optimal if $F^*(i, \lambda) = F^\pi(i, \lambda)$ for every $(i, \lambda) \in S_B$.

We define the following sets of functions: let \mathcal{F} be the set of functions F from S_B into B such that $F(i, \cdot)$ is measurable on B for each $i \in S$, $F(\cdot, \lambda)$ is bounded for each $\lambda \in B$ and $F(i, \lambda) \geq \lambda$ for each $(i, \lambda) \in S_B$, and let \mathcal{F}_ℓ be the set of functions $F \in \mathcal{F}$ such that $F(i, \cdot)$ is nondecreasing and left continuous on B for each $i \in S$. In Theorem 3.1 it is shown that $F^* \in \mathcal{F}_\ell$. However, it is not necessarily true that $F^\pi \in \mathcal{F}_\ell$ for each $\pi \in C$.

We finally define operators T^a , T^δ and T from \mathcal{F} into itself as follows. For $F \in \mathcal{F}$, $(i, \lambda) \in S_B$, $a \in A$ and $\delta \in \Delta_M$,

$$\begin{aligned} T^a F(i, \lambda) &= \sum_{j \in S} \sum_{y \in E} F(j, \lambda \circ y) p^a(j, y|i), \\ T^\delta F(i, \lambda) &= \sum_{a \in A} T^a F(i, \lambda) \delta(a|i, \lambda), \\ TF(i, \lambda) &= \inf_{\delta \in \Delta} T^\delta F(i, \lambda) = \min_{a \in A} T^a F(i, \lambda). \end{aligned}$$

We also define operators T^n by $T^1 = T$ and $T^{n+1} = T(T^n)$, $n \geq 1$. Similarly, $(T^\delta)^n$ is defined for $\delta \in \Delta_M$. In all argument, for $F, G \in \mathcal{F}$, $F \geq G$ means that $F(i, \lambda) \geq G(i, \lambda)$ for all $(i, \lambda) \in S_B$.

3. Optimal value and optimal policy

In this section we prove that the optimal value function is a unique solution to an optimality equation and we give a value iteration method. These results are an associative extension of Eaton and Zadeh[3], Derman[4, 5], and Bellman and Zadeh[1], and a stochastic one of Maruyama[?]. We also show that there exists an optimal left continuous policy.

We first give fundamental lemmas below.

Lemma 3.1.

- (i) For $F, G \in \mathcal{F}$ and $\delta \in \Delta$, $T^\delta F - T^\delta G = T^\delta(F - G)$.
- (ii) If $F, G \in \mathcal{F}$ and $F \geq G$, then $T^a F \geq T^a G$ for each $a \in A$, $T^\delta F \geq T^\delta G$ for each $\delta \in \Delta$ and $TF \geq TG$.
- (iii) If $G \in \mathcal{F}_\ell$, then $T^a G \in \mathcal{F}_\ell$ for any $a \in A$. Also, T is an operator from \mathcal{F} (or \mathcal{F}_ℓ) into itself.
- (iv) If $G_n \in \mathcal{F}_\ell$ and $G_n \leq G_{n+1}$ for each $n \geq 0$, then $\lim_{n \rightarrow \infty} G_n \in \mathcal{F}_\ell$.

Proof. The statements (i) and (ii) are immediate results of definitions.

(iii) Let $G \in \mathcal{F}$. Since $G(i, \cdot)$ is measurable on B , $T^a G(i, \cdot)$ is also measurable for each $a \in A$ and so is $TG(i, \cdot)$. Also, it is obvious that $TG(i, \cdot)$ is bounded for each $\lambda \in B$. Thus $TG \in \mathcal{F}$.

Next, let $G \in \mathcal{F}_\ell$ and let $i \in S$ be arbitrary. Then it easily follows from the definition that $T^a G(i, \lambda) \geq \lambda$ for each $a \in A$ and hence $TG(i, \lambda) \geq \lambda$ for every $(i, \lambda) \in S_B$, since $G(i, \lambda) \geq \lambda$ for every $(i, \lambda) \in S_B$. Also, it follows that $T^a G(i, \cdot)$ is nondecreasing on B for $a \in A$ and hence so is $TG(i, \cdot)$, since $G(i, \cdot)$ is nondecreasing on B . Also, we see by the dominated convergence theorem that $T^a G(i, \cdot)$ is left continuous on B for each $a \in A$, since $G(i, \cdot)$ is left continuous. Thus since A is finite, $TG(i, \cdot)$ is also left continuous on B . Therefore, we have $T^a G \in \mathcal{F}_\ell$ for each $a \in A$ and $TG \in \mathcal{F}_\ell$.

(iv) It is clear that $\lim_n G_n(i, \lambda)$, say $G(i, \lambda)$, is nondecreasing in λ for each $i \in S$ and $G(i, \lambda) \geq \lambda$ for every $(i, \lambda) \in S_B$. Hence we need to establish that $G(i, \cdot)$ is left continuous for each $i \in S$. Let $(i, \lambda) \in S_B$ be arbitrarily fixed and let $\varepsilon > 0$ be arbitrary. Since $G(i, \cdot) = \lim_n G_n(i, \cdot)$, there is an integer \hat{N} such that $G(i, \lambda) - G_n(i, \lambda) < \varepsilon/2$ for every $n \geq \hat{N}$. Also, since $G_n(i, \cdot)$ is left continuous, we see, for each $n \geq \hat{N}$, that there is $\hat{\delta} > 0$ such that $G_n(i, \lambda) - G_n(i, \lambda') < \varepsilon/2$ when $0 < \lambda - \lambda' < \hat{\delta}$ and $\lambda' \in B$. Thus we have

$$\begin{aligned} 0 &\leq G(i, \lambda) - G(i, \lambda') \\ &= G(i, \lambda) - G_n(i, \lambda) + G_n(i, \lambda) - G_n(i, \lambda') + G_n(i, \lambda') - G(i, \lambda') \\ &< \varepsilon \end{aligned}$$

since $G_n(i, \lambda') \leq G(i, \lambda')$. Hence $G(i, \cdot)$ is left continuous. \square

Remark. From the proof of Lemma 3.1 (iv) we notice that Lemma 3.1 (d) in [12] mistakes. Thus in the lemma, $\lim_n G_n \in \mathcal{F}_r$ should be $\lim_n G_n \in \mathcal{F}_\ell$. Furthermore, through the paper [12], \mathcal{F}_r should be \mathcal{F}_ℓ and “right continuous” should be “left continuous”.

We easily see that for each $F \in \mathcal{F}$, there is a measurable decision rule $\delta \in \Delta_D$ satisfying $TF = T^\delta F$, since TF is measurable and A is finite.

Furthermore, the following lemma is important for main theorems.

Lemma 3.2. *For each $F \in \mathcal{F}_\ell$, there exists a left continuous decision rule $\delta \in \Delta_D$ satisfying $TF = T^\delta F$.*

Proof. Let $F \in \mathcal{F}_\ell$ and $(i, \lambda) \in S_B$ be arbitrarily fixed. From Lemma 3.1, $T^a F(i, \cdot)$ is left continuous on B for each $a \in A$. Since A is finite, we see that there exist $\mu > 0$ and $a \in A$ such that $TF(i, u) = T^a F(i, u)$ for all u satisfying $\lambda - \mu < u \leq \lambda$ and $u \in B$. For such an action a , if we define $\delta \in \Delta_D$ by $\delta(i, u) = a$ for every u so that $\lambda - \mu < u \leq \lambda$ and $u \in B$, then δ is left continuous and $TF(i, \lambda) = T^\delta F(i, \lambda)$. \square

For any $\pi = (\delta_n, n \geq 0) \in C$ and a given history $(i, \lambda, a) \in S_B \times A$, the cut-head policy of π to (i, λ, a) is defined by ${}^1\pi^{(i, \lambda, a)} = (\delta_n^{(i, \lambda, a)}, n \geq 0)$ where $\delta_n^{(i, \lambda, a)}(\cdot | h_n) = \delta_{n+1}(\cdot | (i, \lambda, a), h_n)$ for every $h_n \in H_n$ and each $n \geq 0$. Then we see that ${}^1\pi^{(i, \lambda, a)} \in C$ for a fixed (i, λ, a) . For the sake of simplicity we use a notation:

$$T^{\delta_0} F^{1\pi}(i, \lambda) = \sum_{a \in A} \delta_0(a | i, \lambda) \sum_{j, y} F^{1\pi^{(i, \lambda, a)}}(j, \lambda \circ y) p^a(j, y | i)$$

for each $\pi = (\delta_n, n \geq 0) \in C$ and $(i, \lambda) \in S_B$.

Lemma 3.3. *Let $\pi = (\delta_n, n \geq 0) \in C$ be arbitrary. For each $n \geq 0$, $F_{n+1}^\pi = T^{\delta_0} F_n^{1\pi}$ and $F^\pi = T^{\delta_0} F^{1\pi}$. Especially, $F^\pi = T^\delta F^\pi$ when $\pi = \delta^\infty \in C_D^s$.*

Proof. It follows by Markov property that for any $\pi = (\delta_n, n \geq 0) \in C$,

$$\begin{aligned} T^{\delta_0} F_n^{1\pi}(i, \lambda) &= \sum_{a \in A} \delta_0(a | i, \lambda) \sum_{j, y} F_n^{1\pi^{(i, \lambda, a)}}(j, \lambda \circ y) p^a(j, y | i) \\ &= \sum_{a \in A} \delta_0(a | i, \lambda) \sum_{j, y} E_{(j, \lambda \circ y)}^{1\pi^{(i, \lambda, a)}}[\lambda \circ y \circ Z_n] p^a(j, y | i) \\ &= E_i^\pi[\lambda \circ Z_{n+1}] = F_{n+1}^\pi(i, \lambda). \end{aligned}$$

Similarly, it is easy to see that $F^\pi = T^{\delta_0} F^{1\pi}$. \square

We next give fundamental properties for optimal value functions of finite and infinite horizon cases.

Theorem 3.1. We have the following:

(i) For each $n \geq 0$, $F_n^* \in \mathcal{F}_\ell$ and $\{F_n^*, n \geq 0\}$ satisfies equations :

$$F_0^*(i, \lambda) = \lambda, \quad (i, \lambda) \in S_B, \quad F_n^* = TF_{n-1}^*, \quad n \geq 1.$$

(ii) For each $n \geq 0$, there exists a left continuous policy $\pi \in C_D$ such that $F_n^* = F_n^\pi$.

(iii) For each $n \geq 0$, $F_n^* \leq F_{n+1}^* \leq \lim_{n \rightarrow \infty} F_n^* \leq F^*$ and $\lim_{n \rightarrow \infty} F_n^* \in \mathcal{F}_\ell$.

Remark. On the statement (iii) we have $\lim_{n \rightarrow \infty} F_n^* = F^*$ under some conditions, which we will prove in Theorem 3.2.

Proof. We prove the statements (i) and (ii) of this lemma by induction. When $n = 0$, from the fact that $Z_0 = Y_0 = e$ we see that $F_0^*(i, \lambda) = \inf_{\pi \in C} E_i^\pi[\lambda \circ e] = \lambda = F_0^\pi(i, \lambda)$ for any left continuous policy $\pi \in C_D$ and every $(i, \lambda) \in S_B$ and hence $F_0^* \in \mathcal{F}_\ell$, which implies that (i) and (ii) hold for $n = 0$. Assume that these statements are true for $n = k$. Thus, $F_k^* \in \mathcal{F}_\ell$ and there exists a left continuous policy $\sigma \in C_D$ such that $F_k^* = F_k^\sigma$. It follows from Lemma 3.2 that there exists a left continuous decision rule $\delta \in \Delta_D$ such that $TF_k^* = T^\delta F_k^*$, which implies that $\pi = (\delta, \sigma)$ is a left continuous policy in C_D . It also follows from Lemma 3.3 that for each (i, λ) ,

$$F_{k+1}^*(i, \lambda) \leq F_{k+1}^\pi(i, \lambda) = T^\delta F_k^\sigma(i, \lambda) = T^\delta F_k^*(i, \lambda) = TF_k^*(i, \lambda).$$

Conversely, we see from Lemma 3.3 again that for any policy $\tau = (\delta_n, n \geq 0) \in C$,

$$F_{k+1}^\sigma(i, \lambda) = T^{\delta_0} F_k^{\tau_1}(i, \lambda) \geq T^{\delta_0} F_k^*(i, \lambda) \geq T F_k^*(i, \lambda).$$

Taking infimum over $\sigma \in C$, we obtain $F_{k+1}^*(i, \lambda) \geq T F_k^*(i, \lambda)$. Thus, combining with the previous inequality, we have $T F_k^* = F_{k+1}^* = F_{k+1}^\pi$. Hence, π satisfies $F_{k+1}^* = F_{k+1}^\pi$, and from Lemma 3.1(iii), we have $F_{k+1}^* \in \mathcal{F}_\ell$. By induction, the proof of the statements (i) and (ii) is complete.

(iii) Since (B, \circ) is nondecreasing, we have $Z_n \leq Z_{n+1} \leq Z$. Thus we obtain that $F_n^\pi \leq F_{n+1}^\pi \leq F^\pi$, which implies that $F_n^* \leq F_{n+1}^* \leq \lim_{n \rightarrow \infty} F_n^* \leq F^*$. Also, since $F_n^* \in \mathcal{F}_\ell$ by (i), and $F_n^* \leq F_{n+1}^*$ for each $n \geq 0$, it follows from Lemma 3.1(iv) that $\lim_{n \rightarrow \infty} F_n^* \in \mathcal{F}_\ell$. \square

From Theorem 3.1, we have $F_n^* = T^n F_0^*$ for each $n \geq 0$. In order to prove that $F^* = \lim_{n \rightarrow \infty} F_n^*$, we need the following important lemma.

Lemma 3.4. *Let $\pi = \delta^\infty \in C_D^s$ be a policy satisfying condition that for each $(i, \lambda) \in S_B$ there is a constant $M > 0$ such that $P_{(i, \lambda)}^\pi(\lambda \circ Z \leq M) = 1$.*

- (i) *Let $F, G \in \mathcal{F}$. If $F - G \leq T^\delta(F - G)$ on $\{K\}^c \times B$ and $F = G$ on $\{K\} \times B$, then $F \leq G$ on S_B .*
- (ii) *F^π is the unique solution in \mathcal{F} to equation $F = T^\delta F$ with $F(K, \lambda) = \lambda$ for every $\lambda \in B$.*

Proof. (i) Since $F = G$ on $\{K\} \times B$ and the state K is absorbing and cost-free, it follows that

$$T^\delta(F - G)(K, \lambda) = (F - G)(K, \lambda) = 0$$

for every $\lambda \in B$. By the fact we also see that if $(i, \lambda) \in \{K\}^c \times B$, then

$$T^\delta(F - G)(i, \lambda) = \sum_{(j, y) \in \{K\}^c \times E} (F - G)(j, \lambda \circ y) p^{\delta(i, \lambda)}(j, y|i)$$

By a similar argument and induction, it easily follows that $(T^\delta)^n(F - G)(K, \lambda) = 0$ for any $\lambda \in B$ and

$$\begin{aligned} (T^\delta)^n(F - G)(i, \lambda) &= \sum_{(i_1, y_1) \in \{K\}^c \times E} \sum_{(i_2, y_2) \in \{K\}^c \times E} \cdots \sum_{(i_n, y_n) \in \{K\}^c \times E} (F - G)(i_n, \lambda \circ y_1 \circ y_2 \cdots \circ y_n) \\ &\quad \times p^{\delta(i, \lambda)}(i_1, y_1|i) p^{\delta(i_1, \lambda \circ y_1)}(i_2, y_2|i_1) \cdots p^{\delta(i_{n-1}, \lambda \circ y_1 \cdots \circ y_{n-1})}(i_n, y_n|i_{n-1}) \end{aligned}$$

for any $(i, \lambda) \in \{K\}^c \times B$. From the condition, it may be follows that $\lambda \circ y_1 \cdots \circ y_n \leq M$, $P_{(i, \lambda)}^\pi$ -a.s. for all $n \geq 1$, since $\lambda \circ Z_n = \lambda \circ Y_1 \cdots \circ Y_n \leq \lambda \circ Z \leq M$, $P_{(i, \lambda)}^\pi$ -a.s.. Thus we see from the boundedness of $F - G$ that for $\lambda \in B$ there is $L = L(\lambda) > 0$ such that $(F - G)(i_n, \lambda \circ y_1 \circ y_2 \cdots \circ y_n) \leq L$ for all $i_n \in S$, all $y_j \in E, j = 1, 2, \dots$ and all $n \geq 1$, $P_{(i, \lambda)}^\pi$ -a.s.. Then it follows that when $(i, \lambda) \in \{K\}^c \times B$

$$\begin{aligned} (T^\delta)^n(F - G)(i, \lambda) &\leq L \sum_{(i_1, y_1) \in \{K\}^c \times E} \cdots \sum_{(i_n, y_n) \in \{K\}^c \times E} p^{\delta(i, \lambda)}(i_1, y_1|i) \cdots \\ &\quad \cdots p^{\delta(i_{n-1}, \lambda \circ y_1 \cdots \circ y_{n-1})}(i_n, y_n|i_{n-1}) \end{aligned}$$

By the way, it follows that

$$\sum_{(i_1, y_1) \in \{K\}^c \times E} p^{\delta(i, \lambda)}(j, y|i) = P_{(i, \lambda)}^\pi(X_1 \in \{K\}^c).$$

For $n \geq 1$, assume that

$$\begin{aligned} &\sum_{(i_1, y_1) \in \{K\}^c \times E} \cdots \sum_{(i_n, y_n) \in \{K\}^c \times E} p^{\delta(i, \lambda)}(i_1, y_1|i) \cdots p^{\delta(i_{n-1}, \lambda \circ y_1 \cdots \circ y_{n-1})}(i_n, y_n|i_{n-1}) \\ &= P_{(i, \lambda)}^\pi\left(\bigcap_{k=1}^n \{X_k \in \{K\}^c\}\right) \end{aligned}$$

for any $(i, \lambda) \in \{K\}^c \times B$. Then, it follows from Markov property that when $(i, \lambda) \in \{K\}^c \times B$

$$\begin{aligned}
& \sum_{(i_1, y_1) \in \{K\}^c \times E} \cdots \sum_{(i_{n+1}, y_{n+1}) \in \{K\}^c \times E} p^{\delta(i, \lambda)}(i_1, y_1 | i) \cdots p^{\delta(i_n, \lambda \circ y_1 \cdots \circ y_n)}(i_{n+1}, y_{n+1} | i_n) \\
&= \sum_{(i_1, y_1) \in \{K\}^c \times E} P_{(i_1, \lambda \circ y_1)}^\pi \left(\bigcap_{k=1}^n \{X_k \in \{K\}^c\} \right) p^{\delta(i, \lambda)}(i_1, y_1 | i) \\
&= P_{(i, \lambda)}^\pi \left(\bigcap_{k=1}^{n+1} \{X_k \in \{K\}^c\} \right)
\end{aligned}$$

Thus, by induction, we have

$$(F - G)(i, \lambda) \leq (T^\delta)^n (F - G)(i, \lambda) \leq LP_{(i, \lambda)}^\pi \left(\bigcap_{k=1}^n \{X_k \in \{K\}^c\} \right),$$

for every $(i, \lambda) \in \{K\}^c \times B$ and all $n \geq 1$. Since $P_{(i, \lambda)}^\pi(X_n = K \text{ for some } n \geq 1) = 1$ from the assumption so that $1, 2, \dots, K-1$ are transient and K is absorbing, we obtain

$$\lim_{n \rightarrow \infty} P_{(i, \lambda)}^\pi \left(\bigcap_{k=1}^n \{X_k \in \{K\}^c\} \right) = 1 - P_{(i, \lambda)}^\pi \left(\bigcup_{k=1}^{\infty} \{X_k = K\} \right) = 0.$$

Letting $n \rightarrow \infty$ on the above inequality, we have $(F - G)(i, \lambda) \leq 0$ for every $(i, \lambda) \in \{K\}^c \times B$, which completes the proof of the statement (i).

(ii) From the condition it follows that $\lambda \leq F^\pi(i, \lambda) \leq M$ for each $(i, \lambda) \in S_B$ and hence $F^\pi \in \mathcal{F}$. Let $F \in \mathcal{F}$ be a solution to $F = T^\delta F$ with $F = \lambda$ on $\{K\} \times B$. Since F^π satisfies $F^\pi = T^\delta F^\pi$ and $F^\pi = \lambda$ on $\{K\} \times B$, we have $F - F^\pi = T^\delta(F - F^\pi)$ on $\{K\}^c \times B$ and $F = F^\pi$ on $\{K\} \times B$. Thus the statement (i) implies that $F = F^\pi$. \square

Now we are in a position to give a main theorem.

Theorem 3.2. *Suppose that there exists at least one policy $\sigma \in C$ such that for each $(i, \lambda) \in S_B$ there is a constant $M > 0$ such that $P_{(i, \lambda)}^\sigma(\lambda \circ Z \leq M) = 1$.*

(i) $F^* = \lim_{n \rightarrow \infty} F_n^*$.

(ii) F^* is the unique solution in \mathcal{F} to $F = TF$ with $F(K, \lambda) = \lambda$ for every $\lambda \in B$.

(iii) There exists a left continuous policy $\pi = \delta^\infty \in C_D^s$ satisfying $F^* = T^\delta F^*$ on $\{K\}^c \times B$ and π is optimal.

Proof. Let $G^* = \lim_{n \rightarrow \infty} F_n^*$, which is in \mathcal{F}_ℓ from Theorem 3.1. Also, we see from the condition that $F^\sigma(i, \lambda) \leq M$ and hence $\lambda \leq G^*(i, \lambda) \leq F^*(i, \lambda) \leq F^\sigma(i, \lambda) \leq M$ for each $(i, \lambda) \in S_B$. Thus we have $G^*, F^* \in \mathcal{F}$. We first prove that G^* is the unique solution to $F = TF$ with $F(K, \lambda) = \lambda$ for every $\lambda \in B$. Since $F_n^*(K, \lambda) = \lambda$ for each $n \geq 0$ and all $\lambda \in B$, we have $G^*(K, \lambda) = \lambda$. Also, since $F_{n+1}^* = TF_n^* \leq TG^*$ from Theorem 3.1, letting $n \rightarrow \infty$ we obtain $G^* \leq TG^*$. To show the reverse inequality, we fix $(i, \lambda) \in S_B$. Since S and E are finite sets, it follows that for each $\varepsilon > 0$ there is an integer L such that $F_n^*(j, \lambda \circ y) > G^*(j, \lambda \circ y) - \varepsilon$ for all $n \geq L$ and every $(j, y) \in S \times E$. Also, for such n we see that there exist $\hat{a} = \hat{a}(i, \lambda) \in A$ such that $TF_n^*(i, \lambda) = T^{\hat{a}}F_n^*(i, \lambda)$. Thus we have

$$\begin{aligned}
G^*(i, \lambda) &\geq F_{n+1}^*(i, \lambda) = TF_n^*(i, \lambda) = T^{\hat{a}}F_n^*(i, \lambda) \\
&= \sum_{j, y} F_n^*(j, \lambda \circ y) p^{\hat{a}}(j, y | i) \\
&> \sum_{j, y} (G^*(j, \lambda \circ y) - \varepsilon) p^{\hat{a}}(j, y | i) \\
&= T^{\hat{a}}G^*(i, \lambda) - \varepsilon \\
&\geq TG^*(i, \lambda) - \varepsilon.
\end{aligned}$$

Letting $\varepsilon \rightarrow 0$, we have $G^* \geq TG^*$, which implies that G^* is a solution to $F = TF$ with $F(K, \lambda) = \lambda$ for every $\lambda \in B$. Next we prove the uniqueness for G^* . From Lemma 3.2, we see that there is a decision rule $\delta \in \Delta_D$ such that $G^* = TG^* = T^\delta G^*$, since $G^* \in \mathcal{F}_\ell$. Let $F \in \mathcal{F}_\ell$ be another solution to equation $F = TF$ with $F(K, \lambda) = \lambda$ for every $\lambda \in B$. It follows from Lemma 3.2 again that there is a decision rule $\delta' \in \Delta_D$ such that $F = T^{\delta'} F$. Thus we have $G^* = T^\delta G^* \leq T^{\delta'} G^*$ and $F = T^{\delta'} F \leq T^\delta F$ on $\{K\}^c \times B$. Hence we see that $F - G^* \leq T^\delta (F - G^*)$, $G^* - F \leq T^{\delta'} (G^* - F)$ on $\{K\}^c \times B$ and $F = G^*$ on $\{K\} \times R$. From Lemma 3.4(i), we thus obtain $F = G^*$ on S_B . Hence G^* is the unique solution in \mathcal{F} to $F = TF$ with $F(K, \lambda) = \lambda$ for every $\lambda \in B$.

Now we show the statements (i), (ii) and (iii). From Lemma 3.2, there is a left continuous decision rule $\delta \in \Delta_D$ such that $G^* = TG^* = T^\delta G^*$. Also, it follows from Lemma 3.4(ii) that F^π is a unique solution to $F = T^\delta F$ with $F(K, \lambda) = \lambda$ for every $\lambda \in B$, where $\pi = \delta^\infty$. Thus the uniqueness of G^* implies that $G^* = F^\pi$. However, $G^* \leq F^* \leq F^\pi$ from Theorem 3.1. Hence we obtain $G^* = F^* = F^\pi$, which implies that F^* is the unique solution in \mathcal{F} to $F = TF$ with $F(K, \lambda) = \lambda$ for every $\lambda \in B$ and that π is optimal. \square

From Theorems 3.1 and 3.2 we see that a value iteration is given by $F^* = \lim_{n \rightarrow \infty} T^n F_0^*$ where $F_0^*(i, \lambda) = \lambda$ for each $(i, \lambda) \in S_B$. We give another value iteration in the following theorem.

Theorem 3.3. *Suppose that there is at least one policy $\sigma \in C$ such that for each $(i, \lambda) \in S_B$ there is a constant $M > 0$ such that $P_{(i, \lambda)}^\sigma(\lambda \circ Z \leq M) = 1$. Let $G \in \mathcal{F}$ be a function satisfying $G \leq F^*$. Then $\{T^n G\}$ converges and $\lim_{n \rightarrow \infty} T^n G = F^*$.*

Proof. Since $G \in \mathcal{F}$, we have $F_0^*(i, \lambda) = \lambda \leq G(i, \lambda)$ for every $(i, \lambda) \in S_B$. Hence $T^n F_0^* \leq T^n G$, which leads the inequality $F^* = \lim_n T^n F_0^* \leq \liminf_n T^n G$. Conversely, since $G \leq F^*$ and $F^* = TF^*$, we have $T^n G \leq T^n F^* = F^*$ and hence $\limsup_n T^n G \leq F^*$. Therefore, combining with the previous inequality, we have $\lim_n T^n G = F^*$. \square

4. Policy iteration method

In this section we consider a policy space iteration procedure in our model as follows:

- (i) Select an initial policy $\pi_0 = (\delta_0)^\infty \in C_D^s$.
- (ii) At step n , assume that we have a policy $\pi_n = (\delta_n)^\infty \in C_D^s$ and solve the equation $F = T^{\delta_n} F$ with $F(K, \lambda) = \lambda$ for every $\lambda \in B$ to give a function $F^{\pi_n} \in \mathcal{F}$.
- (iii) If $T^{\delta_n} F^{\pi_n} = TF^{\pi_n}$, stop the procedure. If $T^{\delta_n} F^{\pi_n} \neq TF^{\pi_n}$, go the next step.
- (iv) Find a new policy $\pi_{n+1} = (\delta_{n+1})^\infty \in C_D^s$ by $T^{\delta_{n+1}} F^{\pi_n} = TF^{\pi_n}$.
- (v) Return to step (ii) replacing n by $n + 1$.

From Lemma 3.4(ii) we can uniquely solve the equations in \mathcal{F} at step (ii) under some conditions. We have the following convergence theorem.

Theorem 4.1. *Suppose that there exists at least one policy $\sigma \in C$ such that for each $(i, \lambda) \in S_B$ there is a constant $M > 0$ such that $P_{(i, \lambda)}^\sigma(\lambda \circ Z \leq M) = 1$.*

- (i) The sequence $\{F^{\pi_n}\}$ is nonincreasing and converges to F^* .
- (ii) If $T^{\delta_n} F^{\pi_n} = TF^{\pi_n}$, then F^{π_n} is the optimal value and $\pi_n = (\delta_n)^\infty \in C_D^s$ is an optimal policy.

Proof. (i) Since $F^{\pi_n} = T^{\delta_n} F^{\pi_n}$ for each $n \geq 0$ by Lemma 3.3, we have

$$\begin{aligned} F^{\pi_n} - F^{\pi_{n+1}} &= T^{\delta_n} F^{\pi_n} - T^{\delta_{n+1}} F^{\pi_{n+1}} \\ &\geq T^{\delta_{n+1}} F^{\pi_n} - T^{\delta_{n+1}} F^{\pi_{n+1}} \\ &= T^{\delta_{n+1}} (F^{\pi_n} - F^{\pi_{n+1}}) \end{aligned}$$

and $F^{\pi_n} = F^{\pi_{n+1}} = \lambda$ on $\{K\} \times B$. From Lemma 3.4(i) we see that $F^{\pi_n} \geq F^{\pi_{n+1}}$ and hence the sequence $\{F^{\pi_n}\}$ is nonincreasing. Thus $\{F^{\pi_n}\}$ tends to a function $\tilde{F} \in \mathcal{F}$. We now show that $F^* = \tilde{F}$. Since

$F^{\pi_n} = T^{\delta_n} F^{\pi_n} \geq T F^{\pi_n}$, by letting $n \rightarrow \infty$, we obtain the inequality $\tilde{F} \geq T\tilde{F}$. Conversely, it follows by the policy procedure that

$$T F^{\pi_n} = T^{\delta_{n+1}} F^{\pi_n} \geq T^{\delta_{n+1}} F^{\pi_{n+1}} = F^{\pi_{n+1}},$$

which yields that $T\tilde{F} \leq \tilde{F}$, by letting $n \rightarrow \infty$. Hence $\tilde{F} = T\tilde{F}$. Therefore the uniqueness of F^* leads $\tilde{F} = F^*$.

(ii) When $T^{\delta_n} F^{\pi_n} = T F^{\pi_n}$, we see by the procedure that $\pi_n = \pi_k$ for every $k \geq n$. Thus it follows from (i) that $F^* = \tilde{F} = F^{\pi_n}$ and hence π_n is optimal. \square

5. Numerical examples

We first consider an example for a maximum case and get optimal value and optimal policy by the policy iteration method.

Example 5.1. Let $x \circ y = \max(x, y)$. Let $S = \{1, 2, 3\}$ be a state space and 3 be a target node. Assume that the state 3 is absorbing and cost-free. Also let $A = \{a_1, a_2\}$ be an action space. We give the probability distributions by

$$\begin{aligned} p^{a_1}(2, 2|1) &= \frac{2}{3}, & p^{a_1}(3, 2|1) &= \frac{1}{3}, \\ p^{a_1}(3, 6|2) &= p^{a_2}(2, 4|1) = 1, \\ p^{a_2}(2, 8|2) &= p^{a_2}(3, 3|2) = \frac{1}{2}. \end{aligned}$$

Then we have $B = [2, 8]$ and $e = 2$. We consider a policy space procedure to give an optimal policy. Let $\pi_0 = (\delta_0)^\infty \in C_D^s$ be an initial policy such that $\delta_0(i, \lambda) = a_1$ for every $(i, \lambda) \in S_B$. Solving the equation $F = T^{\delta_0} F$ with $F(3, \lambda) = \lambda$ for every $\lambda \in B$, we have

$$F^{\pi_0}(2, \lambda) = \begin{cases} 6 & (2 \leq \lambda \leq 6) \\ \lambda & (6 < \lambda \leq 8) \end{cases}, \quad F^{\pi_0}(1, \lambda) = \begin{cases} \frac{1}{3}\lambda + 4 & (2 \leq \lambda \leq 6) \\ \lambda & (6 < \lambda \leq 8) \end{cases}$$

We now see that $T^{\delta_0} F^{\pi_0} \neq T F^{\pi_0} = \min(T^{a_1} F^{\pi_0}, T^{a_2} F^{\pi_0})$, since

$$T^{a_1} F^{\pi_0}(2, \lambda) = F^{\pi_0}(2, \lambda), \quad T^{a_2} F^{\pi_0}(2, \lambda) = \begin{cases} \frac{11}{2} & (2 \leq \lambda \leq 3) \\ \frac{\lambda}{2} + 4 & (3 < \lambda \leq 8) \end{cases}.$$

Next, using $T^{\delta_1} F^{\pi_0} = T F^{\pi_0}$, we give a policy $\pi_1 = (\delta_1)^\infty \in C_D^s$ by

$$\begin{aligned} \delta_1(3, \lambda) &= a_1 \\ \delta_1(2, \lambda) &= \begin{cases} a_2 & (2 \leq \lambda \leq 4) \\ a_1 & (4 < \lambda \leq 8) \end{cases}, \\ \delta_1(1, \lambda) &= a_1. \end{aligned}$$

By solving $F = T^{\delta_1} F$ with $F(3, \lambda) = \lambda$, F^{π_1} is given by

$$F^{\pi_1}(2, \lambda) = \begin{cases} \frac{11}{2} & (2 \leq \lambda \leq 3) \\ \frac{\lambda}{2} + 4 & (3 < \lambda \leq 4) \\ 6 & (4 < \lambda \leq 6) \\ \lambda & (6 < \lambda \leq 8) \end{cases}, \quad F^{\pi_1}(1, \lambda) = \begin{cases} \frac{1}{3}\lambda + \frac{11}{3} & (2 \leq \lambda \leq 3) \\ \frac{2}{3}\lambda + \frac{11}{3} & (3 < \lambda \leq 4) \\ \frac{1}{3}\lambda + 4 & (4 < \lambda \leq 6) \\ \lambda & (6 < \lambda \leq 8) \end{cases}$$

We can easily check that $T^{\delta_1} F^{\pi_1}(i, \lambda) = T F^{\pi_1}(i, \lambda)$ for every $(i, \lambda) \in S_B$. Thus we stop the procedure. From Theorem 4.1 we obtain the optimal value $F^* = F^{\pi_1}$ and an optimal policy $\pi_1 = (\delta_1)^\infty$. Therefore, since $e = 2$, we have optimal value in the original problem as follows:

$$F^*(1, 2) = \frac{13}{3}, \quad F^*(2, 2) = \frac{11}{2}, \quad F^*(3, 2) = 2.$$

We next consider an example for a multiplicative case.

Example 5.2. Let $x \circ y = xy$. Let $S = \{1, 2, 3\}$ be a state space and 3 be a target node. Assume that the state 3 is absorbing and cost-free. Also let $A = \{a_1, a_2\}$ be an action space. We give the probability

distributions by

$$\begin{aligned} p^{a_1}(2, 2|1) &= \frac{2}{3}, & p^{a_1}(3, 2|1) &= \frac{1}{3}, \\ p^{a_1}(3, 6|2) &= p^{a_2}(2, 4|1) = 1, \\ p^{a_2}(2, 5|2) &= \frac{1}{16}, & p^{a_2}(3, 3|2) &= \frac{15}{16}. \end{aligned}$$

Then we have $B = [1, \infty)$ and $e = 1$. We consider a policy space procedure to give an optimal value and an optimal policy. Let $\pi_0 = (\delta_0)^\infty \in C_D^s$ be an initial policy such that $\delta_0(i, \lambda) = a_1$ for every $(i, \lambda) \in S_B$. Solving the equation $F = T^{\delta_0} F$ with $F(3, \lambda) = \lambda$ for every $\lambda \in B$, we have

$$F^{\pi_0}(2, \lambda) = 6\lambda, \quad F^{\pi_0}(1, \lambda) = \frac{26}{3}\lambda$$

We now see that $T^{\delta_0} F^{\pi_0} \neq T F^{\pi_0}$, since

$$T^{a_1} F^{\pi_0}(2, \lambda) = F^{\pi_0}(2, \lambda) = 6\lambda, \quad T^{a_2} F^{\pi_0}(2, \lambda) = \frac{75}{16}\lambda$$

Next, using $T^{\delta_1} F^{\pi_0} = T F^{\pi_0}$, we give a policy $\pi_1 = (\delta_1)^\infty \in C_D^s$ by

$$\delta_1(3, \lambda) = a_1, \quad \delta_1(2, \lambda) = a_2, \quad \delta_1(1, \lambda) = a_1.$$

By solving $F = T^{\delta_1} F$ with $F(3, \lambda) = \lambda$, F^{π_1} is given by

$$F^{\pi_1}(2, \lambda) = \frac{45}{11}\lambda, \quad F^{\pi_1}(1, \lambda) = \frac{112}{33}\lambda.$$

We can easily check that $T^{\delta_1} F^{\pi_1}(i, \lambda) = T F^{\pi_1}(i, \lambda)$ for every $(i, \lambda) \in S_B$. Thus we stop the procedure. We obtain the optimal value $F^* = F^{\pi_1}$ and an optimal policy $\pi_1 = (\delta_1)^\infty$. Therefore, since $e = 1$, we have optimal value in the original problem as follows:

$$F^*(1, 1) = \frac{112}{33}, \quad F^*(2, 1) = \frac{45}{11}, \quad F^*(3, 1) = 1.$$

References

- [1] R.E. Bellman, L.A. Zadeh, Decision-making in a fuzzy environment. *Management Science*, 17(1970) B141-B164.
- [2] D.P. Bertsekas, J.N. Tsitsiklis, An analysis of stochastic shortest path problems. *Math. Oper. Res.* 16(1991) 580-595.
- [3] J.H. Eaton, L.A. Zadeh, Optimal pursuit strategies in discrete-state probabilistic systems. *Trans. ASME Ser. D, J. Basic Eng.* 84(1962) 23-29.
- [4] C. Derman, On sequential decisions and Markov chains. *Manage. Sci.* 9(1962) 16-24.
- [5] C. Derman, *Finite state Markovian decision processes*. Academic Press, New York, 1970.
- [6] S. Iwamoto S, T. Fujita, Stochastic decision-making in a fuzzy environment. *J. Operations Research Society of Japan*, 38(1995) 467-482.
- [7] S. Iwamoto S, K. Tsurusaki K, T. Fujita, Conditional decision-making in fuzzy environment. *J. Operations Research Society of Japan*, 42(1999) 198-218.
- [8] S. Iwamoto, K. Tsurusaki, T. Fujita, On Markov policies for minimax decision processes. *J. Math. Anal. Appl.* 253(2001) 58-78.
- [9] Y. Maruyama, Associative shortest and longest path problems. *Bulletin of Informatics and Cybernetics*, 31(1999) 147-163.

- [10] Y. Maruyama, An invariant imbedding approach to associative shortest path problems. *Math. Japonica*, 50(1999) 469-480.
- [11] Y. Ohtsubo, K. Toyonaga, Optimal policy for minimizing risk models in Markov decision processes. *J. Math. Anal. Appl.* 271(2002) 66-81.
- [12] Y. Ohtsubo, Minimizing risk models in stochastic shortest path problems. *Math. Methods Oper. Res.* 57(2003) 79-88.
- [13] Y. Ohtsubo, K. Toyonaga, Equivalence classes for minimizing risk models in Markov decision processes. *Math. Methods Oper. Res.* 60(2004) 239-250.
- [14] Y. Ohtsubo, Optimal threshold probability in undiscounted Markov decision processes with a target set. *Applied Math. Computation*, 149(2004) 519-532.
- [15] Y. Ohtsubo, Multistage Markov decision processes with minimum criteria of random rewards. *Bulletin of Informatics and Cybernetics*, 38(2006) 15-25.
- [16] N.G.F. Sancho, Routing problems and Markovian decision processes. *J. Math. Anal. Appl.* 105(1985) 76-83.
- [17] D.J. White, Markov decision processes. John Wiley, New York, 1993.
- [18] D.J. White, Minimizing a threshold probability in discounted Markov decision processes. *J. Math. Anal. Appl.* 173(1993) 634-646.
- [19] C. Wu, Y. Lin, Minimizing risk models in Markov decision processes with policies depending on target values. *J. Math. Anal. Appl.* 231(1999) 47-67.
- [20] H.J. Zimmermann, Fuzzy set theory- and its applications. Kluwer Academic Publishers, Boston, 1996.