

## 制約付マルコフ決定過程：ベクトル最適化によるアプローチ

長岡工業高等専門学校 浦田和芳 (Kazuyoshi WAKUTA)

### 1. はじめに

制約を持つ割引コスト型マルコフ決定過程を考える. 2種類のコスト  $c^0$  と  $c^1$  がある. 各初期状態ごとに期待合計割引  $c^1$  コストについて制約が与えられて, すべての初期状態について期待合計割引  $c^0$  コストを最小にする. 研究の多くは平均コスト型についてであるが, 割引コスト型についてもいくつかの研究がある. Kallenberg[5], Altman と Shwartz[1], Altman[2]のLP法, Frid[4], Sennott[7]のラグランジュ乗数法, Liu と Liu[6]のベクトル最適化法などである. Liu と Liu[6]は, 制約付マルコフ決定過程をベクトル値マルコフ決定過程と関連付け, 確定的定常政策の中から制約最適な政策を求める政策改良法を示した. しかし, 確率的でシステムの履歴に依存するすべての政策が使用可能ならば, 制約最適な政策は確定的定常政策とは限らない. この場合, Liu と Liu[6]のアルゴリズムは適用されない.

本論では, Liu と Liu[6]と同様に制約付マルコフ決定過程をベクトル値マルコフ決定過程と関連付ける. しかし, 確率的でシステムの履歴に依存するすべての政策を考え, 制約最適な政策の存在とその求め方について議論する.

### 2. 制約付マルコフ決定過程

制約付マルコフ決定過程(CMDP)

$S = \{0, 1, 2, \dots, N-1\}$  : 状態空間,  $A(i)$  = 有限集合 : 行動空間

$p(j|i, a), i, j \in S, a \in A(i)$  : 推移確率

$c^0(i, a), c^1(i, a)$  : コスト関数

$\beta$  ( $0 \leq \beta < 1$ ) : 割引因子.

$\Pi$  : すべての政策の集合,  $\Pi_S$  : すべての確率的定常政策の集合,  $\Pi_D$  : すべての確定的定常政策の集合, とおく.

$$I_\pi^0(i_0) = E_\pi \left[ \sum_{n=0}^{\infty} \beta^n c^0(i_n, a_n) | i_0 \right], \quad I_\pi^1(i_0) = E_\pi \left[ \sum_{n=0}^{\infty} \beta^n c^1(i_n, a_n) | i_0 \right], \quad i_0 \in S.$$

$d = (d_0, \dots, d_{N-1})$  : 制約ベクトル.

$$\Delta_{i_0} = \{ \pi \in \Pi \mid I_\pi^1(i_0) \leq d_{i_0} \}, \quad \Delta = \bigcap_{i_0 \in S} \Delta_{i_0}.$$

$I_\pi^0(i_0) \leq I_\pi^1(i_0)$ ,  $\pi \in \Delta_{i_0}$  ならば,  $\pi^*$  は  $i_0$ -制約最適であるという. すべての  $i_0 \in S$  について  $i_0$ -制約最適であるとき,  $\pi^*$  は制約最適であるという.

### 3. ベクトル値マルコフ決定過程との関連

$U \subset R^2$  に対して,  $e(U) = \{x \in U \mid \text{ある } y \in U \text{ に対して } y \leq x \text{ ならば } y = x\}$  とおく.  
 $c(i, a) = (c^0(i, a), c^1(i, a))$  をコスト関数にもつベクトル値マルコフ決定過程(VMDP)を考える.

$$I_\pi(i_0) = E_\pi \left[ \sum_{n=0}^{\infty} \beta^n c(i_n, a_n) \mid i_0 \right], \quad i_0 \in S.$$

$$V(i_0) = \bigcup_{\pi \in \Pi} \{I_\pi(i_0)\}, \quad i_0 \in S, \quad V_D(i_0) = \bigcup_{f \in \Pi_D} \{I_f(i_0)\}, \quad i_0 \in S.$$

このとき,  $V(i_0) = \text{co} V_D(i_0)$ ,  $i_0 \in S$ . すべての  $i_0 \in S$  に対して  $I_{\pi^*}(i_0) \in e(V(i_0))$  であるとき,  $\pi^*$  はVMDPで最適であるという.

VMDPに対して,  $c^\lambda(i, a) = \langle \lambda, c(i, a) \rangle$ ,  $\lambda \in R^2$  をコスト関数にもつスカラー値マルコフ決定過程(MDP( $\lambda$ ))を考える.

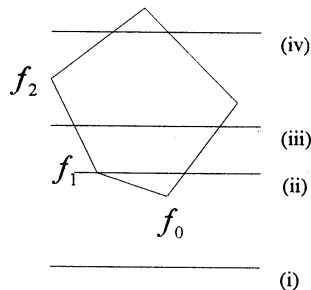
$$I_\pi^\lambda(i_0) = E_\pi \left[ \sum_{n=0}^{\infty} \beta^n c^\lambda(i_n, a_n) \mid i_0 \right], \quad i_0 \in S.$$

$$\pi^* : \text{MDP}(\lambda) \text{で最適} \Leftrightarrow I_{\pi^*}^\lambda(i_0) \leq I_\pi^\lambda(i_0), \quad \forall i_0 \in S, \forall \pi \in \Pi.$$

**Theorem 3.1.**  $E$  をVMDPで最適なすべての確定的定常政策の集合とし,  $E$  の政策を  $I_f^1(i_0)$  の大きさの順に並べる:  $I_{f_0}^1(i_0) \leq I_{f_1}^1(i_0) \leq \dots \leq I_{f_r}^1(i_0)$ . このとき

- (i)  $I_{f_0}^1(i_0) > d_{i_0}$  ならば,  $\Delta_{i_0} = \emptyset$ .
- (ii)  $I_{f_k}^1(i_0) = d_{i_0}$  ならば,  $f_k$  は  $i_0$ -制約最適である.
- (iii)  $I_{f_k}^1(i_0) < d_{i_0} < I_{f_{k+1}}^1(i_0)$  ならば,  $i_0$ -制約最適な(確率的)定常政策が存在する.
- (iv)  $I_{f_r}^1(i_0) \leq d_{i_0}$  ならば,  $f_r$  は  $i_0$ -制約最適である.

**Proof.** (i) (ii) (iv)は明らか.



**Lemma 3.1.**  $\{\lambda_n\} \rightarrow \lambda$  ( $\lambda_n, \lambda \in R^2$ ) で,  $f$  は  $\text{MDP}(\lambda_n)$ ,  $\forall n$  について最適とする. このとき,  $f$  は  $\text{MDP}(\lambda)$  で最適である.

**Lemma 3.2.**  $e_0, e_1$  は,  $V(i_0)$  の 1 つの辺上の有効端点とする. このとき,  $e_0, e_1$  に対応し, 共通の  $\text{MDP}(\lambda)$ ,  $\lambda > 0$  で最適な  $g_0, g_1$  が存在する.

$f, g \in \Pi_D$ ,  $t (0 \leq t \leq 1)$  に対して  $\pi = (t, f, g)$  を, 確率  $t$  で  $f$ , 確率  $1-t$  で  $g$  をとる政策とする.

**Lemma 3.3.**  $f, g$  は  $\text{MDP}(\lambda)$  で最適で,  $I_f^1(i_0) < d_{i_0} < I_g^1(i_0)$  とする. このとき,  $\pi^* = (t^*, f, g)$  が  $\text{MDP}(\lambda)$  で最適で,  $I_{\pi^*}^1(i_0) = d_{i_0}$  となる  $t^* (0 \leq t^* \leq 1)$  が存在する.

**Lemma 3.4.** Lemma 3.3 の条件を仮定する.  $f$  と  $g$  は 1 つの状態だけで異なると仮定する. このとき, Lemma 3.3 の  $t^*$  は一意に定まる.

**Theorem 3.1 (iii) の Proof.**

$e_0 = (e_0^0, e_0^1), e_1 = (e_1^0, e_1^1)$  を  $V(i_0)$  の 1 つの辺上の有効端点とすると  $e_0^1 < d_{i_0} < e_1^1$ . Lemma 3.2 より,  $e_0, e_1$  に対応し  $\text{MDP}(\lambda)$ ,  $\lambda > 0$  で最適な  $g_0, g_1 \in \Pi_D$  が存在する. ここで,

$$I_{g_0}^1(i_0) < d_{i_0} < I_{g_1}^1(i_0).$$

$$h_l \in \Pi_D, l = 0, 1, \dots, N$$

$$h_l(i) = \begin{cases} g_0(i), & i \geq l \\ g_1(i), & i < l. \end{cases}$$

とおく. Chitgopeker [3] より,  $h_l, l = 0, 1, \dots, N$  は  $\text{MDP}(\lambda)$ ,  $\lambda > 0$  で最適で,  $I_{h_l}^1(i_0) < d_{i_0} < I_{h_{l+1}}^1(i_0)$  なる  $l$  が存在する. Lemma 3.4 より,  $\pi^* = (t^*, h_l, h_{l+1})$  が  $\text{MDP}(\lambda)$ ,  $\lambda > 0$  で最適で,  $I_{\pi^*}^1(i_0) = d_{i_0}$  となる  $t^*$  が存在する.  $I_{\pi^*}^0(i_0)$  は制約を満たす政策の中で最小の利得なので,  $\pi^* = (t^*, h_l, h_{l+1})$  は  $i_0$ -制約最適である.

**Remark 3.1.**  $E$  は,  $\text{VMDP}$  における政策反復によって定まる. また, Theorem 3.1 (iii) の  $i_0$ -制約最適な (確率的) 定常政策は, 実際に求めることが可能である. したがって, 制約最適な政策が存在すれば, それを求めることができる.

## 4. 数値例

$$S = \{0,1,2\}, A(0) = A(1) = \{0,1\}, A(2) = \{0\}$$

$$p(0|0,0) = 1, p(1|0,1) = 1$$

$$p(1|1,0) = 1, p(2|1,1) = 1$$

$$p(2|2,0) = 1$$

$$c(0,0) = (0,0), c(0,1) = (-3,0)$$

$$c(1,0) = (2,2), c(1,1) = (4,1)$$

$$c(2,0) = (4,1)$$

$$\beta = 0.5.$$

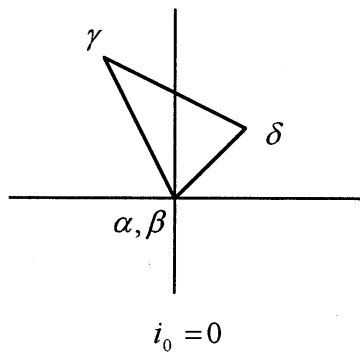
$$\alpha: \alpha(0) = 0, \alpha(1) = 0, \alpha(2) = 0; \beta: \beta(0) = 0, \beta(1) = 1, \beta(2) = 0$$

$$\gamma: \gamma(0) = 1, \gamma(1) = 0, \gamma(2) = 0; \delta: \delta(0) = 1, \delta(1) = 1, \delta(2) = 0.$$

$E = \{\alpha, \beta, \gamma\}$  : 最適な確定的定常政策

$$I_\alpha(0) = (0,0), I_\alpha(1) = (4,4), I_\alpha(2) = (8,2); I_\beta(0) = (0,0), I_\beta(1) = (8,2), I_\beta(2) = (8,2)$$

$$I_\gamma(0) = (-1,2), I_\gamma(1) = (4,4), I_\gamma(2) = (8,2).$$



$E$  の政策を  $I_f^1(0)$  大きさの順に並べる.

$$0 = I_\alpha^1(0) = I_\beta^1(0) < I_\gamma^1(0) = 2.$$

- (i)  $d_0 < 0$  ならば,  $\Delta_0 = \phi$ .
- (ii)  $d_0 = 0$  ならば,  $\alpha$  と  $\beta$  が 0-制約最適.
- (iv)  $d_0 \geq 2$  ならば,  $\gamma$  が 0-制約最適.
- (iii)  $0 < d_0 < 2$  のとき

$\langle \lambda, I_\gamma(0) - I_\alpha(0) \rangle = 0$ となる $\lambda > 0$ を求める.

$I_\gamma(0) - I_\alpha(0) = (-1, 2)$ なので $\lambda = (2, 1)$ .

$E^\lambda$ を求める.

$\langle \lambda, I_\alpha(0) \rangle = \langle \lambda, I_\beta(0) \rangle = \langle \lambda, I_\gamma(0) \rangle = 0$

$\langle \lambda, I_\alpha(1) \rangle = 12, \langle \lambda, I_\beta(1) \rangle = 18, \langle \lambda, I_\gamma(1) \rangle = 12$

$\langle \lambda, I_\alpha(2) \rangle = \langle \lambda, I_\beta(2) \rangle = \langle \lambda, I_\gamma(2) \rangle = 18.$

$E^\lambda = \{\alpha, \gamma\}.$

$E^\lambda$ の中から1つの状態だけで異なる政策のペアを求める.

$\alpha$ と $\gamma$ は1つの状態のみで異なる.

0-制約最適な $\pi = (t, \alpha, \gamma)$ を求める.

$$I_\pi^1(i) = c^1(i, \alpha(i)) + \frac{1}{2} \sum_{j \in S} p(j|i, \alpha(i)) I_\pi^1(j), \quad i = 1, 2; \quad (1)$$

$$I_\pi^1(0) = t c^1(0, \alpha(0)) + (1-t) c^1(0, \gamma(0)) + \frac{1}{2} \sum_{j \in S} (t p(j|0, \alpha(0)) + (1-t) p(j|0, \gamma(0))) I_\pi^1(j). \quad (2)$$

$I_\pi^1(0) = d$ とにおいて, (1)を解くと,  $I_\pi^1(1) = 4, I_\pi^1(2) = 2$ . (2)に代入して

$$d = \frac{1}{2} t d + \frac{1}{2} (1-t) \cdot 4. \quad \therefore t = (4 - 2d) / (4 - d).$$

## 参考文献

- [1] E. Altman and A. Shwartz, Sensitivity of constrained Markov decision processes, Ann. Oper. Res. 32(1991) 1-22.
- [2] E. Altman, Denumerable constrained Markov decision processes and finite approximations, Math. Oper. Res., 19(1994) 169-191
- [3] S.S. Chitgopekar, Denumerable state Markovian sequential control processes : On randomizations of optimal policies, Naval Res. Logistic. Quart. 22(1975) 567-573.

- [4] B.Frid, On optimal strategies in control problems with constraints, *Theory Probab. Appl.* 17(1972)188-192.
- [5] L.C.M.Kallenberg, *Linear Programming and Finite Markovian Control Problems*, Mathematical Centre Tracts No. 148 (Mathematisch Centrum, Amsterdam, 1983).
- [6] J. Liu and K. Liu, Markov decision programming with constraints, *Acta Math. Appl. Sinica*, 10(1994) 1-11.
- [7] L.I.Sennott, Constrained discounted Markov decision chains, *Probab. Engineer. Infor. Sci.*, 5(1991) 463-475