

Functional Characterization for Average Cost Markov

Decision Processes with Doeblin's Conditions

千葉大・教育 蔵野正美 (Masumi Kurano)

Kurano [5,6] は Doeblin's Condition の考え方 [4] を用いて、Compact Metric Space 上の平均コスト基準のマールコフ決定過程 (MDPs) を取扱い、任意の randomized stationary Policy によって誘導される (induced) マールコフ過程に、幾つかの ergodic classes & transient set が許されると general (multichain) Case に対して、最適定常政策の存在定理を与えている。

本報告では、同じ問題設定のもとでの平均コスト基準のマールコフ決定過程に対する階級的特徴 (最適方程式の導出、及び、その有効性) を行なう。なお、この報告は、Kurano [7] の内容を一部手直したものである。

## 1. 定式化

ある complete separable な距離空間のボレル部分集合を準

：ボレル集合とよぶ。ボレル集合  $X$  のボレル部分集合の全体を  $\mathcal{B}_X$  で表す。

マールコフ決定過程 (Markov Decision Processes, MDPs)  
は次の4つの要素  $S, A, Q, C$  から成る：

- (i)  $S$  はボレル集合で状態空間を表す。
- (ii) 各  $x \in S$  に対して,  $A(x)$  はボレル集合  $A$  の部分集合で、状態  $x$  においてとりうる行動の全体を表す。
- (iii)  $C: S \times A \rightarrow (-\infty, \infty)$  は、有界なボレル可測関数で、直接費用関数 (immediate cost function) を表す。
- (iv)  $Q$  は  $\mathcal{B}_S \times S \times A$  上の確率核で次の条件 (a), (b) を満たす。
  - (a) 各  $(x, a) \in S \times A$  に対して,  $Q(\cdot | x, a)$  は  $\mathcal{B}_S$  上の確率測度である。
  - (b) 各  $D \in \mathcal{B}_S$  に対して,  $Q(D | \cdot)$  は  $S \times A$  上のボレル可測関数である。

この報告を通じて、次が仮定される。

### 仮定

- (i)  $S \in R := \{(x, a) \mid x \in S, a \in A(x)\}$  は共にコンパクト集合である。
- (ii) エスト関数  $C$  は非負値有界で、下半連続である。
- (iii)  $x_m \rightarrow x, a_m \rightarrow a$  のとき,  $Q(\cdot | x_m, a_m)$  は  $Q(\cdot | x, a)$  に弱収束する。

考察する決定過程の標本空間は  $\Omega = (\mathcal{S} \times A)^\infty$  で、 $t$  期の状態と行動は、確率変数  $X_t, \Delta_t$  で表す。 $(t \geq 0)$

各  $t \geq 0$  に対して、 $\mathcal{B}_{A \times \mathcal{S} \times (A \times \mathcal{S})^t}$  上の確率核  $\pi_t$  で

$$\pi_t(A(x_t) | x_0, a_0, \dots, a_{t-1}, x_t) = 1$$

for all  $(x_0, a_0, \dots, a_{t-1}, x_t) \in \mathcal{S} \times (A \times \mathcal{S})^t$

を満たす  $x_t$  の集合  $\pi = (\pi_0, \pi_1, \dots)$  を政策といふ。

今、 $D \in \mathcal{B}_S$  に対して、 $T(A|D) \in \mathcal{B}_{A \times D}$  上の確率核  $\pi$  で、  
すべての  $x \in D$  で、 $\pi(A(\cdot) | x) = 1$  を満たす  $\pi$  の全体とする。

政策  $\pi = (\pi_0, \pi_1, \dots)$  がランダム定常政策であるとは、  
ある  $\pi \in T(A|S)$  が存在して、すべての  $(x_0, a_0, \dots, x_t) \in \mathcal{S} \times (A \times \mathcal{S})^t$  とすべての  $t \geq 0$  に対して、

$$\pi_t(\cdot | x_0, a_0, \dots, x_t) = \pi(\cdot | x_t)$$

が成り立つときをいふ。この場合、 $\pi$  を単純  $\pi^{(0)}$  で表す。

任意の  $D \in \mathcal{B}_S$  に対して、 $u(x) \in A(x)$ ,  $x \in D$  を満たすボレル可測 (analytically measurable) 関数  $u: D \rightarrow A$  の全体  $\in \mathcal{B}(D \rightarrow A)$  ( $B_a(D \rightarrow A)$ ) で表す。

ランダム定常政策  $\pi^{(0)}$  が定常であるとは、 $f \in \mathcal{B}(S \rightarrow A)$   
が存在して、 $\pi(f(x)) = 1$ ,  $x \in S$  が成り立つときをいふ。

そのような政策を  $f^{(0)}$  で表す。

$t$  期までの履歴を  $H_t = (X_0, \Delta_0, \dots, \Delta_{t-1}, X_t)$  とする。

任意に与えられた政策  $\pi = (\pi_0, \pi_1, \dots)$  に対して、次の後定

すみ： すべての  $D_1 \in \mathcal{B}_A$ ,  $D_2 \in \mathcal{B}_S$  に対して

$$\text{Prob}(\Delta_t \in D_1 | H_t) = \pi_t(D_1 | H_t)$$

$$\text{Prob}(X_{t+1} \in D_2 | H_{t+1}, \Delta_{t+1}, X_t = x, \Delta_t = a) = Q(D_2 | x, a) \\ (t \geq 0).$$

このとき、任意の政策  $\pi \in \Pi$  と初期分布  $v \in P(S)$  に対して、  
 $\Omega$  上の確率測度  $P_\pi^v$  が通常の方法で定義される。

但し、 $D \in \mathcal{B}_S$  に対して、 $P(D)$  は  $D$  上の確率分布の全体を表す。

次の平均コスト基準を考察する。

任意の  $\pi \in \Pi$  と初期分布  $v \in P(S)$  に対して、

$$\gamma(v, \pi) := \limsup_{T \rightarrow \infty} E_\pi^v [\sum_{t=0}^{T-1} c(X_t, \Delta_t)] / T.$$

但し、 $E_\pi^v$  は  $P_\pi^v$  に関する期待値を表す。

之から、次を定義する。

$$\gamma(v) := \inf_{\pi \in \Pi} \gamma(v, \pi),$$

$$\gamma^* := \inf_{v \in P(S)} \gamma(v).$$

任意の  $D \in \mathcal{B}_S$  に対して、

$$\gamma(x, \pi^*) \leq \gamma(x, \pi), \quad x \in D, \quad \pi \in \Pi$$

が成り立つ。よって、 $\pi^*$  は  $D$  における最適 (Optimal in  $D$ ) である。

$S$  における最適な政策は常に最適である。

## 2 最適方程式 (1)

この節では、MDPs に対する positive recurrence の条件のもとで、最適方程式を導出し、その有効性を議論する。

任意の重  $\pi \in T(A|S)$  に対して、 $t$  期の指移確率  $Q^{(t)}$  を次で定義する：

$$Q^{(0)}(\cdot|x, \pi) = \int Q(\cdot|x, a) \pi(da|x)$$

$$Q^{(t+1)}(\cdot|x, \pi) = \int Q^{(t)}(\cdot|x_t, \pi) Q^{(t)}(dx_t|x, \pi) \quad (t \geq 1).$$

この報告を通じて、次の Doeblin 条件が成り立つことを仮定する。

### 仮定 (Doeblin [4])

次を満足する  $\mathcal{B}_S$  上の有限測度  $\gamma$  と  $\varepsilon > 0$  が存在する：

任意の重  $\pi \in T(A|S)$  に対して、自然数  $l$  が存在して、

$\gamma(D) \leq \varepsilon$  なら  $D \in \mathcal{B}_S$  に対して、 $Q^{(l)}(D|x, \pi) \leq 1 - \varepsilon$  がすべての  $x \in S$  に対して成り立つ。

次の定理はすでに証明されている。

### 定理 2.1 ([5])

次の (i), (ii) を満たす  $\gamma(C) > \varepsilon$  なら  $C \in \mathcal{B}_S$  と定常政策  $\bar{f}^{(0)}$  が存在する。

(i)  $\bar{f}^{(0)}$  は  $C$  において最適で、かつ、すべての  $x \in C$  において  $\gamma^* = \gamma(x, \bar{f}^{(0)})$ 。

(ii) すべての  $x \in C$  において、 $Q(C|x, \bar{f}(x)) = 1$  で、  
 $Q(\cdot|x, \bar{f}(x))$  はまた  $C$  上で誘導された (induced)  $C$  上のスムーズ

フ過程は transient state を持たない。

定理 2.1 の  $C \in \mathcal{B}_S$  に対して、次の Lemma 2.1 が成り立つ。

### Lemma 2.1

定理 2.1 の  $C \in \mathcal{B}_S$  と定常政策  $\pi^{(\infty)}$  に対して、 $C$  上の一様有界なボレル可測関数  $u \in B(C)$  が存在して、次を満たす：

$$u(x) + \gamma^* = C(x, f(x)) + \int_C u(x') Q(dx' | x, f(x))$$

for all  $x \in C$ .

(略証)

Doeblin 条件のもとで  $u$  はフ過程の性質を利用してある。

$$Q(\cdot | x) := Q(\cdot | x, f(x))$$

$\{Q(\cdot | x), x \in C\}$  は  $\exists$  induced Markov Process on  $C$  は one ergodic set のみであると仮定してよい。

$C_1, C_2, \dots, C_d$  : the cyclically moving classes in  $C$

$x \in C_1$  に対して、

$$\lim_{t \rightarrow \infty} Q^{(td+j-1)}(\cdot | x) = v_j(\cdot)$$

uniformly and exponentially fast

$$v(\cdot) := \frac{1}{d} \sum_{j=1}^d v_j(\cdot)$$

$$z \in z \in, \quad \gamma^* = \int_C C(x, f(x)) v(dx) \quad \text{が成り立つ}.$$

$$U^T(x) := \sum_{t=0}^{[T/a]d-1} E_{f^{(ta)}}^x [C(X_t, \Delta_t) - \gamma^*]$$

$$U(x) := \lim_{T \rightarrow \infty} U^T(x) \quad (\text{一様収束})$$

となること、

$$U(x) = \sum_{t=0}^{\infty} E_{f^{(n)}}^x [ C(X_t, A_t) - \gamma^t ]$$

上式を再帰式に書きかければ、  $U(x)$  は Lemma 2.1 の関係式を満たす。  
(証終)

任意の  $D \in \mathcal{B}_S$  に対する hitting time を次で定義する。

$$T_D := \inf \{ t \geq 0 \mid X_t \in D \} \quad \text{但し } \inf \emptyset = \infty$$

Lemma 2.1 の結果を  $C$  が全空間に拡大するためには、  
MDP に対する次の positive recurrence の条件を必要とする。

### 仮定 A

$\gamma(D) > \varepsilon$  なら 任意の  $D \in \mathcal{B}_S$  と  $x \in S$  に対して、

$E_{\pi}^x(T_D) < \infty$  なら  $\pi \in \Pi$  が存在する。

### Lemma 2.2 ([9])

定常政策  $f^{(\infty)}$  と  $D \in \mathcal{B}_S$  に対して、 次の不等式を満たす

$\varphi \in B_c(S)$ ,  $\varphi \geq 0$  と正の定数  $\alpha$  が存在するとき：

$$\varphi(x) \geq \alpha + \int_{S-D} \varphi(x') Q(dx' | x, f(\omega)) , \quad \forall x \notin D$$

となる  $E_{f^{(\infty)}}^x(T_D) \leq \varphi(x)/\alpha$  成立。

### Lemma 2.3

仮定 A のもとで、  $\gamma(D) > \varepsilon$  なら 任意の  $D \in \mathcal{B}_S$  に対して、

$E_{f^{(\infty)}}^x(T_D) < \infty$ ,  $x \notin D$  を満たす  $\bar{f}^{(\infty)} \in B_a(S \rightarrow A)$  が存在する。

ボレル集合  $X$  に対して、 $X$  上の universally measurable functions の全体を  $B_u(X)$  と表す。記述を簡単にするために、 $B_u(S)$  上の operator  $\cup$  を次のようにも定義する：

各  $x \in S$ ,  $a \in A(x)$  に対して

$$\cup(x, a, u) := c(x, a) + \int u(x') Q(dx' | x, a), \quad u \in B_u(S)$$

次の結果はよく知られてる。

#### Lemma 2.4 (Ergodicity [8])

定常政策  $f^{(x)}$  に対して、次の 2 つの条件式を満たす  $u \in B_u(S)$  が存在する：

$$u(x) + \gamma^* \geq \cup(x, f(x), u), \quad \forall x \in S$$

$$\lim_{T \rightarrow \infty} E_{f^{(x)}}^x (u(X_T)) / T = 0$$

このとき、 $f^{(x)}$  は最適政策である。

以上の Lemmas を用いて次の定理を得る。証明は [7] を参照のこと。

定理 2.2 仮定 A のもとで、次の (i) ~ (iii) を満たす  $v \in B_u(S)$  が存在する：

(i)  $v$  は次の最適不等式を満たす。

$$v(x) + \gamma^* \geq \inf_{a \in A(x)} \cup(x, a, v), \quad \forall x \in S$$

(ii) 定常政策  $f^{(x)}$  ( $f \in B_a(S \rightarrow A)$ ) が次の (\*), (\*\*\*) を満たすならば、 $f^{(x)}$  は最適である。

$$(**) \quad v(x) + \gamma^* \geq \cup(x, f(x), v) \quad \forall x \in S$$

$$(***) \quad \lim_{T \rightarrow \infty} E_{f(x)}^x (v(X_T)) / T = 0$$

(iii) (\*\*) (\*\*\*) を満たす定常政策  $f^{(n)}$  ( $f \in B_a(S \rightarrow A)$ ) が存在する  
3.

### 3. 最適方程式 (2)

この節は positive recurrence の仮定 (前節の仮定 A) が成り立たない場合について考察する。

#### 定義

任意の  $D \in \mathcal{B}_S$  に対して,

$$\gamma(D) := \left\{ x \in S - D \mid E_\pi^x(T_D) < \infty \text{ for some } \pi \in \Pi \right\}$$

#### 定義

任意の  $D \in \mathcal{B}_S$  に対して.

$$P(D) := \left\{ (v, \pi) \in P(D) \times \Pi \mid P_\pi^v(X_t \in D \text{ for all } t \geq 0) \right\}$$

$$\gamma^*(D) := \inf_{(v, \pi) \in P(D)} \gamma(v, \pi)$$

但し  $P(D) = \emptyset$  のとき  $\gamma^*(D) = \infty$  とする。

次の仮定が必要である:

仮定 B 次の B1, B2 が成り立つ。

B1.  $\gamma(D) > 0$  かつ 3 任意の  $D \in \mathcal{B}_S$  に対して.

$Q(\partial D \mid x, a) = 0, \forall x \in S, a \in A(x)$  但し  $\partial D$  は  $D$  の boundary を表す。

B2. 任意の  $D \in \mathcal{B}_S$  に対して.  $Q(D \mid x, a)$  は  $(x, a) \in S \times A$  の連続関数。

次の Lemma の証明は [6] の Lemma 3.4 やもしくは、前節の定理 2.2 と Lemma 2.3 に基づいています。

### Lemma 3.1

$P(G) \neq \emptyset$  かつ  $G \in \mathcal{B}_S$  に対して、次の (i) ~ (iii) を満たす定常政策  $\bar{f}^{(\infty)}$ 、 $f(C) > 0$  かつ  $C \in \mathcal{B}_G$  もしくは、 $v \in B_u(Y(C) \cup C)$  が存在する：

$$(i) \quad \gamma(x, \bar{f}^{(\infty)}) = \gamma^*(G) \quad \forall x \in Y(C) \cup C.$$

$$(ii) \quad Q(Y(C) \cup C | x, \bar{f}^{(\infty)}) = 1, \quad \forall x \in Y(C) \cup C.$$

(iii)  $v$  は  $C \cup \bar{C}$  一様有界で次の 2 つの連続式を満たす：

$$v(\alpha) + \gamma^*(G) = c(\alpha, \bar{f}^{(\infty)}) + \int_{Y(C) \cup C} v(\alpha') Q(d\alpha' | x, \bar{f}^{(\infty)})$$

for all  $x \in Y(C) \cup C$

$$\lim_{T \rightarrow \infty} E_{\bar{f}^{(\infty)}}^x (v(X_T)) / T = 0.$$

### 定義

任意の  $D \in \mathcal{B}_S$  に対して、 $A(x, D) := \{a \in A^\infty \mid Q(D | x, a) = 1\}$

以上、準備の上とて次の定理を得る。証明は [7] を参照のこと。

定理 3.1 仮定 B の下で、次の (i) (ii) を満たす  $S$  の可測分割と  $v \in B_u(S)$  が存在する：

$$S = S_1 \cup S_2 \cup \dots \cup S_r \cup F, \quad F \in \mathcal{B}_S, \quad S_i \in \mathcal{B}_S, \quad S_i \cap S_j = \emptyset \text{ (i)}$$

(i) 各  $i$  ( $1 \leq i \leq r$ ) に対して、次の最適不等式が成立する。

$$V(x) + \gamma^*(S_i^*) \geq \inf_{a \in A(x, S_i)} U(x, a, v) \quad x \in S_i, \bigcup_{i=1}^r S_i^* = \bigcup_{i=1}^r S_i.$$

(ii)  $T \in \mathbb{Z}_{\geq 1}$ , 次の最適方程式が成り立つ。

$$V(x) = \inf_{a \in A(x)} \left\{ \sum_{j=1}^r \gamma^*(S_j^*) Q(S_j | x, a) + \int V(x') Q(dx' | x, a) \right\}$$

#### Reference.

- [1] Bertsekas, D.P. and Shreve, S.D. (1978). Stochastic Optimal Control - The Discrete Time Case, Academic Press.
- [2] Borkar, V.S. (1983). Controlled Markov chains and stochastic networks. Siam J. Control and Optimization. 21 652-666.
- [3] ————. (1984). On minimum cost per unit time control of Markov chains. Siam J. Control and Optimization. 22 965-978.
- [4] Doob, J.L. (1953). Stochastic Processes, Wiley, New York.
- [5] Kurano, M. (1989). The existence of a minimum pair of state and policy for Markov decision processes under the hypothesis of Doeblin. Siam J. Control and Optimization. 27 296-307.
- [6] ————. (1989) Average cost Markov decision processes under the hypothesis of Doeblin. Technical Reports of Mathematical Sciences, Chiba University. Vol.4 (1988) No.9.  
To appear in Annals of Operations Research.
- [7] ————. (1990). Functional characterization for average cost Markov decision processes with Doeblin's conditions. Technical Reports of Mathematical Sciences, Chiba University, Vol.6 (1990) No.1. To appear in Comp.& Math. Appl..
- [8] Ross, S.M. (1970). Applied Probability Models with Optimization Applications. Holden-Say, San Francisco.
- [9] Tweedie, R.L. (1976). Criteria for classifying general Markov chains. Adv. Appl. Prob., 8, 737-771.