

マルコフ型決定過程(II)

九大理 古川長太

§1. 序

Wald以来のsequential analysisとBellmanのDynamic Programmingの中のdiscrete parameterの部分を、まとめたMarkovian Decision Processとして最初に定式化したのはBlackwellである。

sequential analysisと所謂D.P.との相異は、前者においてはstopping ruleとstopしたときのdecisionの二通り、一方が之よりsequential decision ruleの追究が問題であるのに反し、後者においては殆んどstopping ruleが問題にされない実に在る。而るに前者は吸收壁をもつマルコフ過程とみなされることから、結局両者は同じ範囲に属する問題として定式化されることがある。

この様を一般的に定式化では、従来研究されて来た様な個々の問題に対する解を、又はその近似解を求めると言った具体的な結果を期待出来ないのは当然であるが、反面それらの

問題に共通の本質的な部分を浮きぼりにして、解明すること
が出来る。

ここでは、これらの事について、主として Blackwell []、
Strauch [] の結果と、併せて Furukawa の拡張(大部分)
を附加えて報告する。

§2. 諸定義

Def. 2.1

X, Y ; ある complete separable metric space の Borel subset

\mathcal{B} ; X における Borel 集合族

\mathcal{B} ; Y における Borel 集合族

$P(X)$; X 上の probability distribution の class

$g(y|x)$; conditional probability distribution (各 $x \in X$ について、

$g(\cdot|x)$ は Y 上の prob. dist., 各 $B \in \mathcal{B}$ について

$g(B|\cdot)$ は X 上の \mathcal{B} 可測関数)

$Q(Y|X)$; $g(y|x)$ の class

Def. 2.2

$M(X) : \begin{cases} X \text{上の 実数値有界可測関数の class (D-case)} \\ X \text{上の non-positive 実数値有界可測関数の class (N-case)} \\ X \text{上の non-negative 実数値有界可測関数の class (P-case)} \end{cases}$

(D, N, P の定義は後述)

Def. 2.3 $M(XY)$; Def. 2.2 りおひて 定義空間を直積空間

XY としたもの

Def. 2.4

$$pu; \quad pu \equiv \int_X u(x) d\mu(x) \quad \text{for } \mu \in P(X), \quad u \in M(X)$$

$$qu; \quad qu(x) \equiv \int_Y u(x, y) d\eta(y|x) \quad \text{for } \eta \in Q(Y|X), \quad u \in M(X)$$

Def. 2.5 以上の諸定義を拡張して

$$P(X_1 X_2 \dots X_n), \quad P(X_1 X_2 \dots),$$

$$Q(X_{n+1} | X_1 X_2 \dots X_n), \quad Q(X_{n+1} X_{n+2} | X_1 X_2 \dots X_n),$$

$$M(X_1 X_2 \dots X_n)$$

Def. 2.6

$\mu \in P(X)$ ならば μ が degenerate $\Leftrightarrow \mu\{\cdot|x\} = 1$ for some $x \in X$

$\eta \in Q(Y|X)$ ならば η が degenerate $\Leftrightarrow \eta(\cdot|x)$ が各 $x \in \cup$ で degenerate

Def. 2.7

S ; state space (non-empty Borel set)

A ; action space (non-empty Borel set)

$S \ni s$; state

$A \ni a$; action

$Q(S|SA) \ni \eta$; transition probability distribution

$M(SAS) \ni r$; reward function

$\beta \equiv (\beta_1, \beta_2, \beta_3, \dots)$, ($0 \leq \beta_i \leq 1$, $i=1, 2, \dots$);

discount factor vector

Def. 2.8

(S, A, g, r, β) ; discrete Dynamic Programming の構成要素

Def. 2.9

$H_n \equiv S A S A \cdots S A S$ (2n-1 factors)

$\pi_n \in Q(A|H_n)$, $n=1, 2, \dots \in \mathcal{L}$

$\pi \equiv (\pi_1, \pi_2, \dots)$; policy

Def. 2.10 $e_\pi \equiv \pi_1 g \pi_2 g \cdots \in Q(ASAS \cdots | S)$

$I_n(\pi, \beta, v) \equiv e_\pi \left[\sum_{j=1}^n \beta_1 \beta_2 \cdots \beta_{j-1} r(s_j, a_j, s_{j+1}) + \beta_1 \beta_2 \cdots \beta_n v \right]$, ($v \in M(S)$)

$I_n(\pi, \beta) \equiv I_n(\pi, \beta, 0)$

$I(\pi, \beta) \equiv e_\pi \sum_{j=1}^{\infty} \beta_1 \beta_2 \cdots \beta_{j-1} r(s_j, a_j, s_{j+1})$; expected total reward

Def. 2.11

random Markov policy; π はより $\pi_n \in Q(A|S)$, $n=1, 2, \dots$

Markov policy; random Markov policy π はより 各 π_n が degenerate

Markov policy $\in (f_1, f_2, \dots)$ と書く。

stationary policy; (f, f, f, \dots)

§ 3. 基本定理

Def. 3.1 $u, v \in M(X)$ は $\exists f (f$

$u \geq v \iff u(x) \geq v(x)$ for all $x \in X$

Lem. 3.1 (Blackwell and Nardzewski [])

$g \in Q(Y|X)$, $S \in \sigma \times \beta_Y$, $g(S_x|x) > 0$ for all $x \in X$

$\Rightarrow (x, g(x)) \in S$ for all $x \in X$ すなはち σ 可測 $g(x)$ が存在する。

Theorem 3.1

$\forall g \in Q(Y|X)$, $\forall u \in M(XY)$, $\forall \epsilon > 0$,

\exists degenerate $f \in Q(Y|X)$

; (i) $fu \geq gu$

(ii) $g(\{y : u(x_0, y) \geq u(x_0, f(x_0)) + \epsilon\} | x_0) = 0$
for all $x_0 \in X$

Def. 3.2

D-case ; $\sum_{j=1}^{\infty} \beta_1 \beta_2 \cdots \beta_j \equiv L$: 收束, $r \in M(SAS)$

N-case ; $\beta_i = 1$ for $i = 1, 2, \dots$, $r \in M(SAS)$

P-case ; $\beta_i = 1$ for $i = 1, 2, \dots$, $r \in M(SAS)$, かつ
R-in policy $\pi = \tau + (\tau \text{ exp. total reward } \pi)$ 有界

D-case と更に、次の二つの場合を合併す。

D-h-case ; $\beta_i = b$ for $i = 1, 2, \dots$

D-n-case ; D-case と D-h case 以外

Lem 3.2

(D-case) $I_n(\pi, \beta, v) \rightarrow I(\pi, \beta)$ for $\forall v \in M(S)$ 收束

(P-case) $I_n(\pi) \uparrow I(\pi)$ 收束

(N-case) $I_n(\pi) \downarrow I(\pi)$ ($-\infty$ は発散のところもある)

Def. 3.3

π^* が (p, ε) -optimal ; $\Pr\{I(\pi^*, \beta) \geq I(\pi, \beta) - \varepsilon\} = 1$ for $\forall \pi$

π^* が ε -optimal ; $I(\pi^*, \beta) \geq I(\pi, \beta) - \varepsilon$ for $\forall \pi$

π^* が optimal ; $I(\pi^*, \beta) \geq I(\pi, \beta)$ for $\forall \pi$

ただし N -case では上の定義に, $I(\pi^*, \beta) > -\infty$ を加える。

§4. D-case

Lem. 4.1 (Blackwell [])

各 $\beta \in P(S)$, 各 $\varepsilon > 0$ に対して (p, ε) -optimal policy が存在。

Def. 4.1 π^* が π を (p, ε) -dominate す ;

$$\Pr\{I(\pi^*, \beta) \geq I(\pi, \beta) - \varepsilon\} = 1$$

Lem. 4.2 (Blackwell [])

各 $\beta \in P(S)$, 各 $\varepsilon > 0$, 各 π に対して, π を (p, ε) -dominate す Markov policy が存在する。

Theorem 4.1 (Blackwell [])

各 $\beta \in P(S)$, 各 $\varepsilon > 0$ に対して (p, ε) -optimal Markov policy が存在する。

Def. 4.2 (Furuikawa)

$$T_{n_j}; T_{n_j} u(s) = \int [r(s, f_m(s), t) + \beta_j u(t)] d\gamma_j(t | s, f_m(s))$$

$\pi = (f_1, f_2, \dots)$ に対して

$$U_j; U_j u(s) = \sup_n T_{n_j} u(s)$$

$T_{n_j} \in (f_m, \beta_j)$ に対する operator と云う。 $(f_m, \beta_j) \sim T_{n_j}$ と書く。

$\sqcup_j \in (\pi = (f_1, f_2, \dots), \beta_j)$ に対する operator と云う。
 $(\pi, \beta_j) \sim \sqcup_j$ と書く。

Theorem 4.2 (Furukawa)

(a) T_{n_j} は単調, 即ち $u \leq v \Rightarrow T_{n_j} u \leq T_{n_j} v$

(b) c が定数なら $T_{n_j}(u+c) = T_{n_j} u + \beta_j c$

(c) Markov $\pi = (f_1, f_2, \dots)$ に対しては

$$T_{11} T_{22} \cdots T_{nn} v = I_n(\pi, \beta, v)$$

Def. 4.3

$\pi = (\pi_1, \pi_2, \dots)$ に対して ${}^n\pi \equiv (\pi_{n+1}, \pi_{n+2}, \dots)$

従って $\pi = (f_1, f_2, \dots)$ に対して ${}^n\pi = (f_{n+1}, f_{n+2}, \dots)$

$\beta = (\beta_1, \beta_2, \dots)$ に対して ${}^n\beta \equiv (\beta_{n+1}, \beta_{n+2}, \dots)$

Theorem 4.2 (d)

$\pi = (f_1, f_2, \dots)$ に対して

$$T_{nn} I({}^n\pi, {}^n\beta) = I({}^{n-1}\pi, {}^{n-1}\beta) \quad \text{for } n=1, 2, \dots$$

Def. 4.4 (Blackwell [])

π は Markov policy と云う

f が π -generated; f は $S \rightarrow A$ の mapping

S の partition $\{S_n\}$ が存在して

$f = f_n$ on S_n for each n .

policy $\hat{\pi}$ が π -generated ; $\hat{\pi} = (g_1, g_2, \dots)$ は Markov.
各 g_n が π -generated

Def. 4.5 $F(\pi)$; π -generated function の class
 $G(\pi)$; π -generated policy の class

Theorem 4.3 (Furukawa)

(a) 各 $\pi = (f_1, f_2, \dots)$, 各 $\varepsilon > 0$, 各 β_j に付して $\hat{f}_j \in F(\pi)$ が
存在し, $(\hat{f}_j, \beta_j) \rightsquigarrow \hat{T}_{jj}$ とすれば
 $\hat{T}_{jj} u \geq U_j u - \varepsilon \quad \text{for } \forall u$

(b) π を 任意の Markov policy とする。各 $\hat{f} \in F(\pi)$ に付して
 $(\hat{f}, \beta_j) \rightsquigarrow \hat{T}_j$ とすれば
 $\hat{T}_j u \leq U_j u \quad \text{for } \forall j, \forall u$

Def. 4.6 (Furukawa)

$(\pi, \beta_j) \rightsquigarrow U_j$ とし, $\lim_{n \rightarrow \infty} U_j U_{j+1} \dots U_n u \equiv U_{j-1}^*$ とき
これを (π, β) に付する limit point とす。記号で
 $(\pi, \beta) \rightsquigarrow U_{j-1}^*$ と書く。特に, $U_0^* = U^*$ と書く。

Theorem 4.4 (Furukawa)

(a) π を 任意の Markov policy とする

$$I(\hat{\pi}, \beta) \leq U^* \quad \text{for } \forall \hat{\pi} \in G(\pi)$$

(b) π を 任意の Markov policy とする

$$\forall \varepsilon > 0, \exists \hat{\pi} \in G(\pi); I(\hat{\pi}, \beta) \geq U^* - \varepsilon$$

- (c) 各 β_j ($j \geq 0$) に対して ε -optimal policy (β_j に依存する)
が“有れば”， β に対する $(1+\varepsilon)$ -optimal Markov policy
が存在する。 ($\varepsilon \geq 0$)
- (d) $(f = a, \beta_j)$ に対する operator を T_{aj} とする。 各 $\varepsilon > 0$,
各 β_j に対して ε -optimal policy (β_j に依存する) が“有れば”
Markov policy $\hat{\pi}$ が存在して $(\hat{\pi}, \beta_j) \rightarrow \hat{u}_j^*$ は \sim では
 $\hat{u}_j^* \rightarrow$ 可測関数 i かつ

$$\hat{u}_{j+1}^* = \sup_{a \in A} T_{aj} \hat{u}_j^* \quad \text{for each } j$$

が成立する。
- (e) 各 j について， $\hat{\pi}^*$ が β_j に対して optimal となるための必要十分条件は

$$I(\hat{\pi}^*, \beta_j) = \sup_{a \in A} T_{aj} I(\hat{\pi}^*, \beta_j) \quad \text{for each } j$$

である。（この条件式を， π^* に関する optimality equation と云う）

[証明]

(a) π を任意の Markov とする。

$$\hat{\pi} \equiv (g_1, g_2, \dots) \in G(\pi), \quad (g_j, \beta_j) \rightarrow \hat{T}_j \text{ とする。}$$

$$\text{Theorem 4.2 (d)} \Rightarrow \hat{T}_{nm} I(\hat{\pi}, \beta) = I(\hat{\pi}, \beta) \quad \text{for each } n$$

$$\begin{aligned} \hat{T}_{m-1, m-1} \hat{T}_{mm} I(\hat{\pi}, \beta) &= \hat{T}_{m-1, m-1} I(\hat{\pi}, \beta) \\ &= I(\hat{\pi}, \beta). \end{aligned}$$

$$\therefore \hat{T}_{11} \hat{T}_{22} \dots \hat{T}_{mm} I(\hat{\pi}, \beta) = I(\hat{\pi}, \beta)$$

$$u_m \equiv I(\hat{\pi}, \beta).$$

$$\therefore \hat{T}_{11} \hat{T}_{22} \cdots \hat{T}_{nn} u_n = I(\hat{\pi}, \beta).$$

$$\|u\| \equiv \sup_s |u(s)|$$

-方 任意の $u \in M(S)$ は

$$\begin{aligned} \|\hat{T}_{11} \hat{T}_{22} \cdots \hat{T}_{nn} u - \hat{T}_{11} \hat{T}_{22} \cdots \hat{T}_{nn} u_n\| &\leq \beta_1 \beta_2 \cdots \beta_n \|u_n - u\| \\ &\leq \beta_1 \beta_2 \cdots \beta_n (L \|r\| + \|u\|). \end{aligned}$$

$$\therefore \|I(\hat{\pi}, \beta) - \hat{T}_{11} \hat{T}_{22} \cdots \hat{T}_{nn} u\| \leq \beta_1 \beta_2 \cdots \beta_n (L \|r\| + \|u\|)$$

$$\therefore \hat{T}_{11} \hat{T}_{22} \cdots \hat{T}_{nn} u \rightarrow I(\hat{\pi}, \beta) \quad \dots \dots \dots (1)$$

各 g_j は $g_j \in F(\pi)$ だから Th. 4.3 (b) より

$$\hat{T}_{jj} u \leq U_j u \quad \text{for each } j.$$

$$\therefore \hat{T}_{j+1, j+1} (\hat{T}_{jj} u) \leq \hat{T}_{j+1, j+1} (U_j u) \leq U_{j+1} U_j u.$$

$$\therefore \hat{T}_{11} \hat{T}_{22} \cdots \hat{T}_{jj} u \leq U_1 U_2 \cdots U_j u \quad \text{for all } j$$

左2は (1) より

$$I(\hat{\pi}, \beta) \leq u^*$$

$$(b) \varepsilon' \equiv \varepsilon / (1+L)$$

Th. 4.3 (a) は より β_j, u は depend して $\hat{f}_j \in F(\pi)$ が $T_2 T_1$

$$(\hat{f}_j, \beta_j) \sim \hat{T}_{jj} \text{ と } z$$

$$\hat{T}_{jj} u \geq U_j u - \varepsilon' \quad \dots \dots \dots (2)$$

$\beta_{j+1}, U_j u - \varepsilon'$ は depend して $\hat{f}_{j+1} \in F(\pi)$ が $T_2 T_1$

$$(\hat{f}_{j+1}, \beta_{j+1}) \sim \hat{T}_{j+1, j+1} \text{ と } z$$

$$\hat{T}_{j+1, j+1} (U_j u - \varepsilon') \geq U_{j+1} (U_j u - \varepsilon') - \varepsilon'$$

(2) は より

$$\hat{T}_{j+1, j+1} (\hat{T}_{jj} u) \geq \hat{T}_{j+1, j+1} (U_j u - \varepsilon') \geq U_{j+1} (U_j u - \varepsilon') - \varepsilon'$$

$$= U_{j+1} U_j u - \beta_{j+1} \varepsilon' - \varepsilon'$$

$$\therefore \hat{T}_{11} \hat{T}_{22} \cdots \hat{T}_{jj} u \geq U_1 U_2 \cdots U_j u - \varepsilon' (1 + \sum_{n=1}^{j-1} \beta_1 \beta_2 \cdots \beta_n)$$

$$\therefore I(\hat{\pi}, \beta) \geq u^* - \varepsilon' (1+L) = u^* - \varepsilon.$$

(c) 各 β_j ($j \geq 0$) は ε -optimal policy があるとする。

$\pi^{*\beta} \in \beta$ は ε -optimal policy とする。

$$\pi^{*\beta^{-1}} = (\pi_{j=1}, \pi_{j=2}, \dots)$$

$$\begin{aligned} \therefore I(\pi^{*\beta^{-1}}, \beta) &= \pi_{j=1} g[r + \beta_j I(\pi^{*\beta^{-1}}, \beta)] \\ &\leq \pi_{j=1} g[r + \beta_j \{I(\pi^{*\beta}, \beta) + \varepsilon\}] \\ &\leq f_j g[r + \beta_j \{I(\pi^{*\beta}, \beta) + \varepsilon\}] \quad (\text{Th. 3.1 は式3}) \\ &= T_{jj} I(\pi^{*\beta}, \beta) + \beta_j \varepsilon. \end{aligned}$$

$$\therefore T_{jj} I(\pi^{*\beta}, \beta) \geq I(\pi^{*\beta^{-1}}, \beta) - \beta_j \varepsilon \quad \dots (3)$$

(3) は各 j について成立するから

$$T_{j=1, j=1} I(\pi^{*\beta^{-1}}, \beta) \geq I(\pi^{*\beta^{-2}}, \beta) - \beta_{j=1} \varepsilon$$

左は degenerate $f_{j=1} \neq 1$ で。

$$\begin{aligned} \therefore T_{j=1, j=1} T_{jj} I(\pi^{*\beta}, \beta) &\geq T_{j=1, j=1} I(\pi^{*\beta^{-1}}, \beta) - \beta_{j=1} \beta_j \varepsilon \\ &\geq I(\pi^{*\beta^{-2}}, \beta) - \beta_{j=1} \varepsilon - \beta_{j=1} \beta_j \varepsilon \end{aligned}$$

$$\therefore T_{11} T_{22} \dots T_{jj} I(\pi^{*\beta}, \beta) \geq I(\pi^*, \beta) - \varepsilon \left(\sum_{n=1}^j \beta_1 \beta_2 \dots \beta_n \right)$$

$\hat{\pi} = (f_1, f_2, \dots)$ とする。上式より

$$I(\hat{\pi}, \beta) \geq I(\pi^*, \beta) - \varepsilon L$$

$\rightarrow \pi^*$ は β に対する ε -optimal である

$$I(\pi^*, \beta) \geq I(\pi, \beta) - \varepsilon \quad \text{for } \pi$$

$$\therefore I(\hat{\pi}, \beta) \geq I(\pi, \beta) - \varepsilon(1+L) \quad \text{for } \pi$$

従って $\hat{\pi}$ は β に対する $(1+L)\varepsilon$ -optimal Markov policy である。

(d) 若く ε , 若く β は ε -optimal policy であると仮定する。

従って, 各 $m = 1, 2, \dots$ は β に対する $\frac{1}{m(1+L)}$ -optimal policy が存在する。

Th. 4.4 (c) より β は β に対する $\frac{1}{m}$ -optimal Markov policy である。

すなはち $\pi^{\beta m} \in G(\hat{\pi})$ である。

$(\hat{\pi}, \beta)$ は $\hat{\pi}_j^*$ である。

Th. 4.4 (a) より

$$I(\pi, \beta) \leq \hat{\pi}_j^* \quad \text{for } \pi \in G(\hat{\pi})$$

$$\therefore I(\pi^{\beta m}, \beta) \leq \hat{\pi}_j^* \quad \text{for } m \dots (4)$$

-> π^* の 定義から

$$I(\pi^*, \beta) \geq I(\pi, \beta) - \frac{1}{m} \quad \text{for } \forall \pi \quad \dots \dots (5)$$

(4), (5) より

$$\hat{u}_j^* \geq I(\pi, \beta) - \frac{1}{m} \quad \text{for } \forall \pi$$

$$\therefore \hat{u}_j^* \geq I(\pi, \beta) \quad \text{for } \forall \pi, \forall j \quad \dots \dots (6)$$

Th. 4.4 (b) により、各 m, j に対して $\tilde{\pi}^{mj} \in G(\hat{\pi})$ かつ $T_{\tilde{\pi}} = T_{\hat{\pi}}$

$$I(\tilde{\pi}^{mj}, \beta) \geq \hat{u}_j^* - \frac{1}{m}$$

$$\therefore \hat{u}_j^* \leq I(\tilde{\pi}^{mj}, \beta) + \frac{1}{m}$$

$$\therefore T_{\tilde{\pi}} \hat{u}_j^* \leq T_{\tilde{\pi}} [I(\tilde{\pi}^{mj}, \beta) + \frac{1}{m}]$$

$$= I((a, \tilde{\pi}^{mj}), \beta) + \frac{\beta}{m}$$

$$\leq \hat{u}_{j-1}^* + \frac{\beta}{m}. \quad ((b) \text{ による})$$

$$\therefore \sup_a T_{\tilde{\pi}} \hat{u}_j^* \leq \hat{u}_{j-1}^*.$$

$$-> \sup_a T_{\tilde{\pi}} \hat{u}_j^* \geq \cup_j \hat{u}_j^* = \hat{u}_{j-1}^*.$$

$$\therefore \sup_a T_{\tilde{\pi}} \hat{u}_j^* = \hat{u}_{j-1}^* \quad \text{for } \forall j.$$

(e) π^* があるて 各 π^* が β に対する optimal と仮定する。

Th. 4.4 (c) により π^* は Markov policy とてよ。

$$\pi^* = (f_1^*, f_2^*, \dots)$$

Th. 4.2 (d) による

$$I(\beta^{-1} \pi^*, \beta)_{S_0} = T_{f_1^*(S_0), j} I(\beta \pi^*, \beta)_{S_0} \quad \text{for } \forall j$$

$$\therefore I(\beta^{-1} \pi^*, \beta)_{S_0} \leq \sup_a T_{aj} I(\beta \pi^*, \beta)_{S_0} \quad \text{for } \forall j$$

これは 各 S_0 で 成立するから

$$I(\beta^{-1} \pi^*, \beta) \leq \sup_a T_{aj} I(\beta \pi^*, \beta) \quad \dots \dots (7)$$

->

$$I(\beta^{-1} \pi^*, \beta) \geq I((a, \beta \pi^*), \beta)$$

$$= T_{aj} I(\beta \pi^*, \beta) \quad \text{for } \forall j$$

$$\therefore I(\beta^{-1} \pi^*, \beta) \geq \sup_a T_{aj} I(\beta \pi^*, \beta) \quad \text{for } \forall j \quad \dots \dots (8)$$

(7), (8) より

$$I(\beta^{-1} \pi^*, \beta) = \sup_a T_{aj} I(\beta \pi^*, \beta) \quad \text{for } \forall j \quad \dots \dots (9)$$

(9) は π の optimality equation である。

次に 逆の証明。

(9) を仮定する。

Th. 4.1 によると (π, β) -optimal Markov $\hat{\pi} = (\hat{f}_1, \hat{f}_2, \dots)$ が存在する。

から

$$\Pr\{I(\hat{\pi}, \beta) \geq I(\pi, \beta) - \varepsilon\} = 1 \quad \text{for } \forall \pi$$

次に 今 π^* は

$$I(\hat{\pi}, \beta)_{S_0} \geq I(\pi, \beta)_{S_0} - \varepsilon \quad \text{for } \forall \pi \quad \cdots \cdots (10)$$

$(\hat{f}_j, \beta_j) \sim \hat{T}_{j,j}$ とするに (9) の 1 節の下では

$$\hat{T}_{j,j} I(\hat{\pi}^*, \beta_j) \leq I(\hat{\pi}^*, \beta_j) \quad \text{for } \forall j$$

$$\therefore \hat{T}_{j-1,j-1} \hat{T}_{j,j} I(\hat{\pi}^*, \beta_j) \leq \hat{T}_{j-1,j-1} I(\hat{\pi}^*, \beta_j)$$

$$\leq I(\hat{\pi}^*, \beta_j) \quad \text{for } \forall j \quad \quad \quad (9) \text{ によると}$$

$$\therefore \hat{T}_{11} \hat{T}_{22} \cdots \hat{T}_{jj} I(\hat{\pi}^*, \beta_j) \leq I(\pi^*, \beta) \quad \text{for } \forall j$$

$$\therefore I(\hat{\pi}, \beta) \leq I(\pi^*, \beta) \quad \cdots \cdots (11)$$

(10), (11) より

$$I(\pi, \beta)_{S_0} \leq I(\hat{\pi}, \beta)_{S_0} + \varepsilon \leq I(\pi^*, \beta)_{S_0} + \varepsilon \quad \text{for } \forall \pi$$

$$\therefore I(\pi, \beta)_{S_0} \leq I(\pi^*, \beta)_{S_0} \quad \text{for } \forall \pi$$

これは各 S_0 で成立するから

$$I(\pi, \beta) \leq I(\pi^*, \beta) \quad \text{for } \forall \pi$$

したがって π^* は β に対する optimal policy。

各 $j=0, 1, \dots$ で π^* が β_j に対して optimal であることを Th. 4.1 で証明

が出来た。

[Theorem 4.4 の証明終了]

Corollary 各 β_j ($j \geq 0$) に対して optimal policy が“あれど”、 β に対する optimal Markov policy が“存在する。

Def. 4.7 (Blackwell)

equivalent action ; $a, b \in A$.

$$r(s, a, \cdot) = r(s, b, \cdot) \Leftrightarrow$$

$$g(\cdot | s, a) = g(\cdot | s, b) \quad \text{on } \exists, a \neq b \text{ は}$$

S_1 における equivalent である。

essentially countable by π ;

$\pi = (f_1, f_2, \dots)$, 各 (s, a) に対して, $f_n(s) \times a$ が S_1

における equivalent である n が存在する $\forall s \in S_1$, A は ess. count. by π である。

essentially finite by π ;

$\pi = (f_1, f_2, \dots)$,

S の partition $\{S_n\}$ が存在して, 各 n について $s \in S_n$

を除く各 (s, a) に対して $f_1(s), f_2(s), \dots, f_n(s)$ の中の少なくとも一个 $\rightarrow \infty$ S_1 における a が equivalent である $\forall s \in S_1$, A は ess. finite by π である。

Theorem 4.5 (Furukawa)

(a) A が ess. count. by some π なら, 各 $\epsilon > 0$ に対して ϵ -optimal Markov policy が存在する。

(b) A が ess. finite by some π なら, optimal Markov policy が存在する。

§ 5. D-h case 及び N-case

§ 4 の諸定理において, homogeneous discount factor を γ とすると, 更に詳しい結果が得られる。

この場合, D-h case では, discount factor を改めて $\beta <$

よし。Def. 3.2 の D-case の定義より、当然 $\beta < 1$ でなければならぬ。

Def. 5.1

$T_m : (f_m, \beta)$ に対応する operator

$\sqcup : (\pi, \beta)$ に対応する operator

Lem. 5.1 \sqcup は Banach space において contraction mapping
(ただし, D-h case のみ)

Def. 5.2 \sqcup の fixed point $\in U^*$ である。
(ただし, D-h case のみ)

Theorem 5.1 (D-h), (Blackwell)

(a) 各 $p \in P(S)$, 各 $\varepsilon > 0$ に対して $\exists (p, \varepsilon)$ -Optimal stationary policy
が存在する。

(b) ε -optimal policy が“あれば”, $\varepsilon/(1-\beta)$ -optimal stationary
policy が存在する。
($\varepsilon \geq 0$)

(c) π^* が optimal $\Leftrightarrow I(\pi^*) = \sup_a T_a I(\pi^*)$

Theorem 5.2 (D-h), (Blackwell)

(a) A が ess. count. by some π ならば, 各 $\varepsilon > 0$ に対して
 ε -optimal stationary policy が存在する。

(b) A が ess. finite by some π ならば, optimal stationary
policy が存在する。

Def. 5.3 π^* が 強義 (p, ε) -optimal ; $P\{I(\pi^*) \geq \sup_{\pi} I(\pi) - \varepsilon\} = 1$

Theorem 5.3 (D-h, N) (Strauch)

各 p , 各 $\varepsilon > 0$ に対して, 強義 (p, ε) -optimal Markov policy

が存在する。ただし、N-case では収束する expected total reward をもつ policy π^* が一意に存在することを仮定する。

(註) 一般の N-case では (ρ, ε) -optimal stationary policy π^* の存在がどうか明らかにされていない。

Theorem 5.4 (N) (Strang)

Optimal policy π^* が存在すれば optimal stationary policy π^* が存在する。

Theorem 5.5 (N) (Strang)

A が ess. finite by some π , かつ, 収束する expected total reward をもつ policy π^* が一意に存在すれば, optimal stationary policy π^* が存在する。

§ 6. Additional results.

Theorem 6.1 (P)

各 π , 各 $\varepsilon > 0$, 各 $p \in P(S)$ に対して, 次の式を満たす semi-Markov policy τ 及び Markov policy ς が存在する。;

$$I(\tau) \geq I(\pi) - \varepsilon, \quad p I(\varsigma) \geq p I(\pi) - \varepsilon$$

Theorem 6.2 (P)

各 $p \in P(S)$, 各 $\varepsilon > 0$ に対して, (ρ, ε) -optimal semi-Markov policy π^* が存在する。

(註) Th. 5.4 の P-case でも成立するかどうか不明らかである。

Theorem 6.3 (D, P, N)

$$v^* \equiv \sup_{\pi} I(\pi) < \infty < \infty,$$

$$v^*(s) = \sup_a T_a v^*(s) \quad \text{for all } s \in S.$$

かつ, D-case では v^* は上式の unique bounded solution である。

より, P-case では v^* は上式の, $s_1 \in S$ に一様に最小の non-negative solution である。

Theorem 6.4 (D, N)

$$I(f, \pi) \geq I(\pi) \Rightarrow I(f^{(\infty)}) \geq I(\pi)$$

Theorem 6.5 (D, N)

σ, τ を policy として, π を次の様に定義する。

$$\pi_n = \begin{cases} \sigma_n & \text{if } (s_1, a_1, \dots, a_{n-1}, s_n) \in B_n \\ \tau_n & \text{if } (s_1, a_1, \dots, a_{n-1}, s_n) \in B_n^c \end{cases}$$

$\pi \in \mathcal{L}$.

$$B_n = \{(s_1, a_1, \dots, a_{n-1}, s_n) \mid u_n > v_n\},$$

$$u_n(s_1, a_1, \dots, a_{n-1}, s_n) = \sum_{j=n}^{\infty} \beta^{j-n} \sigma_j g_j \cdots \sigma_n g_n,$$

$$v_n(s_1, a_1, \dots, a_{n-1}, s_n) = \sum_{j=n}^{\infty} \beta^{j-n} \tau_j g_j \cdots \tau_n g_n,$$

$$\Rightarrow I(\pi) \geq \max(I(\sigma), I(\tau))$$

Theorem 6.6 (D, N)

$\sigma = (f_1, f_2, \dots)$, $\tau = (g_1, g_2, \dots)$ を Markov policy とする。

$\pi = (h_1, h_2, \dots)$ を次の様に定める。

$$h_n = \begin{cases} f_n & \text{if } I(\pi^{n-1}\sigma) > I(\pi^{n-1}\tau) \\ g_n & \text{if otherwise} \end{cases}$$

$$\Rightarrow I(\pi) \geq \max(I(\sigma), I(\tau))$$

Theorem 6.7 (D, N)

$f^{(\infty)}, g^{(\infty)}$ は stationary policy として,

$$h = \begin{cases} f & \text{if } I(f^{(\infty)}) \geq I(g^{(\infty)}) \\ g & \text{if otherwise} \end{cases}$$

$$\Rightarrow I(h^{(\infty)}) \geq \max(I(f^{(\infty)}), I(g^{(\infty)}))$$

§ 7. Remarks.

マルエフ型決定過程として定式化される問題としては,
discrete time D.P. problem の他に, Bayesian の統計的
諸問題がある。以下、これらを列挙する。

- 1) Sequential Analysis
- 2) Bayesian Adaptive Control
- 3) Replacement Problem
- 4) Taxicab Problem
- 5) Allocation Problem
- 6) Inventory Problem
- 7) Smoothing Problem