

Optimal Stochastic Control

北大 理 古川 長太

§1. 離散時間の場合

1.1 一般的な定義

X, Y ; ある完備可分距離空間の Borel 部分集合

$P(X)$; X 上のすべての確率測度の族

$P(y|x)$; x を与えたときの Y 上の条件付確率測度

$Q(Y|X)$; 上記の $P(y|x)$ の族

$M(X)$; X 上の実数値有界 Baire 関数の族

以上の諸定義は, 定義空間を拡張して

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

$$Q(Y_1 Y_2 \cdots Y_m | X_1 X_2 \cdots X_n),$$

$$Q(Y_1 Y_2 \cdots | X_1 X_2 \cdots X_n)$$

の形に拡張される。

$$pu; \quad pu = \int_X u(x) dP(x) \quad \text{for } P \in P(X), \quad u \in M(X)$$

$$g u; g u(x) = \int_Y u(x, y) d g(y|x) \quad \text{for } g \in Q(Y|X), u \in M(XY)$$

この定義も、直積空間を定義空間として拡張される。

$p \in P(X)$ なる p が degenerate ; $p\{x\} = 1$ for some $x \in X$

$g \in Q(Y|X)$ なる g が degenerate ; $g(\cdot|x)$ が 各 x に対して degenerate

1.2 最適化問題の定義

S ; 状態空間 (ある完備可分距離空間の Borel 部分集合)

A ; 行動空間 (ある完備可分距離空間の Borel 部分集合)

$S \ni s$: 状態 (state)

$A \ni a$: 行動 (action)

$H_n \equiv S A S A \cdots S A S \quad (2n-1 \text{ 回})$

$g^N \equiv \{g_1, g_2, \dots, g_N\}$ 長 N の $g_n \in Q(S|H_n A) \quad (1 \leq n \leq N)$

$\pi^N \equiv \{\pi_1, \pi_2, \dots, \pi_N\}$ 長 N の $\pi_n \in Q(A|H_n) \quad (1 \leq n \leq N)$

$(1 \leq N \leq +\infty)$

このような π^N を 政策 (policy) とする。

マルコフ政策 (Markov policy) ; $\pi^N = \{\pi_1, \pi_2, \dots, \pi_N\}$ にお

いて, 各 $\pi_n \in Q(A|S)$ であつて, 各 π_n が degenerate

マルコフ政策を $\pi^N = \{f_1, f_2, \dots, f_N\}$ と書く。

定常政策 (stationary policy) ; マルコフ政策において

$$f_1 = f_2 = \dots = f_N$$

$r_m \in M(H_{n+1})$ ($n=1, 2, \dots, N$) ; 利益関数 (reward)

評価関数としてつぎのものを考える。

$$I_N(\pi^N, r^N) = \sum_{n=1}^N \pi_n r_n$$

又は, $p \in P(S)$ として

$$p I_N(\pi^N, r^N).$$

政策 π^N を採用したとき System の運動は (π^N, r^N) で決定される H_{N+1} 上の確率過程になる。

Δ を政策のある与えられた族として, Δ の中で

$$I_N(\pi^N, r^N) \rightarrow \text{Max (又は Min)}$$

$$p I_N(\pi^N, r^N) \rightarrow \text{Max (又は Min)}$$

$$I_\infty(\pi^\infty, r^\infty) \rightarrow C \text{ (given constant)}$$

$$p I_\infty(\pi^\infty, r^\infty) \rightarrow C \text{ (given constant)}$$

ならしめる問題を離散時間確率的系における (one person) 最適化問題という。

1.3 マルコフ型の場合

$$r_1 = r_2 = \dots = r_n = \dots = r$$

$$r \in Q(S|SA)$$

$$V \in M(SAS) \quad (\text{reward})$$

$$\beta \quad (0 < \beta < 1) \quad (\text{discount factor})$$

$$I(\pi) \equiv \sum_{n=1}^{\infty} \beta^{n-1} \pi_1 r \pi_2 r \dots \pi_n r \in M(S)$$

$$\pi^* \text{ が } (p, \varepsilon)\text{-最適} ; \quad p \{ I(\pi^*) \geq \sup_{\pi} I(\pi) - \varepsilon \} = 1$$

$$\pi^* \text{ が } \varepsilon\text{-最適} ; \quad I(\pi^*) \geq \sup_{\pi} I(\pi) - \varepsilon$$

$$\pi^* \text{ が } \text{最適} ; \quad I(\pi^*) \geq \sup_{\pi} I(\pi)$$

定理 1.1 (Strauch [7], Furukawa [2])

各 $p \in P(S)$, 各 $\varepsilon > 0$ に対して, (p, ε) -最適政策が存在する。

(この定理は, 最初 Strauch が証明し, 後に Furukawa が証明と簡単化した)

定理 1.2 (Blackwell [1])

各 $p \in P(S)$, 各 $\varepsilon > 0$ に対して (p, ε) -最適マルコフ政策が存在する。

定理 1.3 (Blackwell [1])

各 $p \in P(S)$, 各 $\varepsilon > 0$ に対して (p, ε) -最適定常政策が

存在する。

定理 1.4 (Blackwell [1])

(i) A が可算集合なら, 各 $p \in P(S)$, 各 $\varepsilon > 0$ に対して ε -最適 定常政策が存在する。

(ii) A が有限集合なら, 各 $p \in P(S)$, 各 $\varepsilon > 0$ に対して 最適定常政策が存在する。

定理 1.5 (Howard [3])

A, S が有限集合なら, 政策反復法によって最適政策の構成が出来る。

1.4 S が吸収点をもつ場合

(Ω, \mathcal{F}, P) 確率空間

$\{\mathcal{F}_n\}_{n=1,2,\dots}$; $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots \subset \mathcal{F}_n \subset \dots$ なる \mathcal{F} の
sub- σ -field の系列

$\{x_n\}_{n=1,2,\dots}$; 各 n につき x_n が \mathcal{F}_n -可測であるような確率過程

$\tau(\omega)$; つぎの (i) ~ (iii) を満たす確率変数

(i) $\tau(\omega)$ は正整数値をとる

(ii) $P\{\tau(\omega) < +\infty\} = 1$

(iii) $\{\omega; \tau(\omega) = n\} \in \mathcal{F}_n$ for each n

このように $\tau(\omega)$ を *stopping time* といい。

$E(x_\tau) \rightarrow \text{Max (又は Min)}$ をらしめる (*stopping time* τ に関して) 問題を *Optimal Stopping Problem* といい。

$X = \{x_n, \mathcal{F}_n\}$, $Y = \{y_n, \mathcal{F}_n\}$ に対して

$$Y < X \quad \text{w.p.1} \stackrel{\text{def}}{\iff} y_n < x_n \quad \text{w.p.1 for all } n$$

maximal semi-Martingale ;

$Z = \{z_n, \mathcal{F}_n\}$ に対して, $X = \{x_n, \mathcal{F}_n\}$ があって

(i) X は *semi-Martingale*

(ii) $X < Z$ w.p.1

(iii) $Y < Z$ w.p.1 なるいかなる *semi-Martingale* Y に

対して, $Y < X$ w.p.1

なるとき, X は Z に関して *maximal semi-Martingale* であるといふ。

定理 1.6 (Snell [6])

$Z = \{z_n, \mathcal{F}_n\}$ をある与えられた確率過程,
 $X = \{x_n, \mathcal{F}_n\}$ を Z に関して *maximal semi-Martingale* と

かつ regular な確率過程とする。 $\tau^*(\omega)$ を

$$\tau^*(\omega) = \begin{cases} 1 & \text{if } x_1(\omega) = \infty \\ j & \text{if } x_{k_1}(\omega) < z_{k_1}(\omega) - \varepsilon \text{ for } k_1 < j \text{ and } x_j(\omega) \geq z_j(\omega) - \varepsilon \\ \infty & \text{その他} \end{cases}$$

を定義すると, $\tau^*(\omega)$ は stopping time τ かつ

$$E(x_\tau) \geq E(x_{\tau^*}) - \varepsilon \quad \text{for all } \tau \in T$$

(ただし T は, すべての stopping time の族)

上のようを考之を control model に適用する。

$$Z^{(i)} = \{z_n^{(i)}, \mathcal{F}_n\}_{n=1,2,\dots} \quad (i=1,2,\dots)$$

ただし 各 $Z^{(i)}$ は Ω 上の実数値確率過程で, $z_n^{(i)}$ は \mathcal{F}_n -可測とする。

$A = \{1, 2, 3, \dots\}$; action space

$a_n(\omega)$; $\Omega \rightarrow A$ への \mathcal{F}_{n-1} -可測な mapping ($n=1,2,\dots$)

$\tau(\omega)$; stopping time

$\pi \equiv \{a_1(\omega), a_2(\omega), \dots, a_{\tau(\omega)}(\omega)\}$; policy

Φ_n ; \mathbb{R}^m 上に定義された \mathcal{B}_n -可測実数値関数 ($n=1,2,\dots$)

δ_{ij} ; Kronecker のデルタ

問題 $Z \equiv \{Z^{(1)}, Z^{(2)}, \dots, Z^{(i)}, \dots\}$ (ただし

$Z^{(i)} = \{z_n^{(i)}\}_{n=1,2,\dots}$) を与えられた確率過程の系列

と L,

$$\Psi_\tau(Z, \pi) \equiv \Phi_\tau \left(\sum_{i_1=1}^{\infty} \delta_{i_1} a_1 z_1^{(i_1)}, \sum_{i_2=1}^{\infty} \delta_{i_2} a_2 z_2^{(i_2)}, \dots, \sum_{i_\tau=1}^{\infty} \delta_{i_\tau} a_\tau z_\tau^{(i_\tau)} \right)$$

とおいて,

$$E[\Psi_\tau(Z, \pi)] \rightarrow \text{Min}$$

から求める (π, τ) に関して 向題を考える。

Def $X^{(i)} = \{x_n^{(i)}\}_{n=1,2,\dots}$ を Ω 上の確率過程と L

$$X \equiv \{X^{(1)}, X^{(2)}, \dots, X^{(i)}, \dots\} \quad \text{とおく。}$$

$\{\Phi_n(x_1^{(i_1)}, x_2^{(i_2)}, \dots, x_n^{(i_n)})\}_{n=1,2,\dots}$ が各 (i_1, i_2, \dots) に対して semi-Martingale であるとき, X は $\{\Phi_n\}$ に関して uniform semi-Martingale であるという。

Def $\tilde{\pi} \equiv \{a_1(\omega), a_2(\omega), \dots\}$

$$\Psi_n(X, \tilde{\pi}) \equiv \Phi_n \left(\sum_{i_1=1}^{\infty} \delta_{i_1} a_1 x_1^{(i_1)}, \sum_{i_2=1}^{\infty} \delta_{i_2} a_2 x_2^{(i_2)}, \dots, \sum_{i_n=1}^{\infty} \delta_{i_n} a_n x_n^{(i_n)} \right)$$

$\{\Psi_n(X, \tilde{\pi})\}$ を $\tilde{\pi}$ によって変換された X 上の process といふ。

定理 1.7 (Furukawa)

$X = \{X^{(1)}, X^{(2)}, \dots\}$ は $\{\Phi_n\}$ に関して uniform semi-Martingale と L,

$$\sum_{i_1=1}^{\infty} \sum_{i_2=1}^{\infty} \cdots \sum_{i_n=1}^{\infty} E [|\Phi_n(x_1^{(i_1)}, x_2^{(i_2)}, \dots, x_n^{(i_n)})|] < \infty$$

for each n ,

$\tilde{\pi}$ によって変換された X 上の process $\{\Psi_n(X, \tilde{\pi})\}$ が P に因して一様可積分とする。

このとき,

$\{\Psi_n(X, \tilde{\pi})\}$ は 4 収束して, かつ, $\{\Psi_n(X, \tilde{\pi}), n \geq 1, \Psi_{\infty}(X, \tilde{\pi})\}$ は semi-Martingale である。

(上に設定した問題, あるいは, もっと拡張された問題に対する解答は次回に発表する予定)

1.5 2-person の最適化問題

S ; state space

A_1 ; player I の action space

A_2 ; player II の action space

$A \equiv A_1 \times A_2$

$Q \in \mathcal{Q}(S | SA)$; transition probability measure

$r \in M(SA)$; reward function

$0 < \beta < 1$; discount factor

$$H_n^{(1)} \equiv SA_1 SA_1 \cdots SA_1 S \quad (2n-1 \text{ 回})$$

$$H_n^{(2)} \equiv SA_2 SA_2 \cdots SA_2 S \quad (2n-1 \text{ 回})$$

$$\pi = \{\pi_1, \pi_2, \dots\} ; \pi_n \in Q(A_1 | H_n^{(1)}) \quad (n=1, 2, \dots)$$

..... player I の policy

$$\sigma = \{\sigma_1, \sigma_2, \dots\} ; \sigma_n \in Q(A_2 | H_n^{(2)}) \quad (n=1, 2, \dots)$$

..... player II の policy

$$\pi \vee \sigma \equiv \{\pi_1, \sigma_1, \pi_2, \sigma_2, \dots\} \quad \text{paired policy}$$

$$I(\pi \vee \sigma) = \sum_{n=1}^{\infty} \beta^{n-1} \pi_1 \sigma_1 \gamma \pi_2 \sigma_2 \gamma \cdots \pi_n \sigma_n \gamma \uparrow$$

player I は $I(\pi \vee \sigma) \rightarrow \text{Max}$ ならぬ (π について)

player II は $I(\pi \vee \sigma) \rightarrow \text{Min}$ ならぬ (σ について)

Def

$\pi^* \vee \sigma^*$ が $(\beta, \epsilon_1, \epsilon_2)$ -最適 ;

$$\beta \{ I(\pi^* \vee \sigma) + \epsilon_1 \geq I(\pi^* \vee \sigma^*) \geq I(\pi \vee \sigma^*) - \epsilon_2 \} = 1 \quad \text{for all } \pi \vee \sigma$$

Def $\{\pi^m\}$ を player I の policy の系列とする。このとき

$$\hat{\pi} \text{ が } \{\pi^m\} \text{ で生成される} \stackrel{\text{def}}{\iff} \exists \{S_m\} \text{ partition of } S ;$$

$$\hat{\pi} = \pi^m \text{ for } s_i \in S_m$$

$\{\pi^m\}$ で生成される policy の全体を $O_I(\{\pi^m\})$ と書く。

同様にして, player II の policy の系列 $\{\sigma^m\}$ で生成される policy の全体 $O_{II}(\{\sigma^m\})$ も定義される。

定理 1.8 (Furukawa)

$$\sup_{\pi} P I(\pi \vee \sigma) = \sup_m P I(\pi^m \vee \sigma) \text{ for every } \sigma \in \mathcal{O}_f(\{\sigma^m\})$$

$$\inf_{\sigma} P I(\pi \vee \sigma) = \inf_m P I(\pi \vee \sigma^m) \text{ for every } \pi \in \mathcal{O}_f(\{\pi^m\})$$

ゆえに、各 $\varepsilon_1 > 0$, 各 $\varepsilon_2 > 0$ に対して $(\rho, \varepsilon_1, \varepsilon_2)$ -最適な paired policy が存在する。

§2 連続時間の場合

2.1 マルコフ型の場合

はじめに Miller ([5]) の結果を紹介する。

$S = \{1, 2, \dots, n\}$; state space

A ; action space である有限集合とする。

$F \equiv A \times A \times \dots \times A$ (n 個)

$\pi(x)$ (n ベクトル) ; $[0, T] \rightarrow F$ なる L -可測な mapping

$f = (f_1, f_2, \dots, f_n) \in F$ に対して

$$r(f) = \begin{pmatrix} r(1, f_1) \\ r(2, f_2) \\ \vdots \\ r(n, f_n) \end{pmatrix}$$

$$Q(f) = \begin{pmatrix} q(i, f_i, j) \\ \vdots \\ q(i, f_i, i) \end{pmatrix} \begin{matrix} \downarrow \\ i \\ \rightarrow j \end{matrix} ; \text{ Markov infinitesimal generator matrix}$$

ただし

$$f(i, a, j) \geq 0 \quad \text{for } i \neq j, \text{ for all } a \in A,$$

$$\sum_{j=1}^n f(i, a, j) = 0.$$

$Q(\pi(t))$ の L -可測性を仮定すれば "殆んどすべての t について (1) を満たす絶対連続一意解 $P(s, t)$ が存在し, かつそれは Markov transition matrix になる。

$$(1) \quad \begin{cases} \frac{d}{dt} P(s, t) = P(s, t) Q(\pi(t)) \\ P(s, s) = I \end{cases}$$

評価関数: $v = \int_0^T P(0, t) \Upsilon(\pi(t)) dt \rightarrow \text{Max}$ ならぬ。

最大原理

policy π が, 補助方程式:

$$(2) \quad \begin{cases} -\frac{d\psi}{dt} = \Upsilon(\pi(t)) + Q(\pi(t))\psi & 0 \leq t \leq T \\ \psi(T) = 0 \end{cases}$$

の一意解 $\psi(t)$ に対して

$$\Upsilon(\pi(t)) + Q(\pi(t))\psi(t) = \max_{f \in F} [\Upsilon(f) + Q(f)\psi(t)], \quad \text{a.e. } t \in [0, T]$$

をみたすとき, π は 最大原理をみたすといふ。

上のことと equivalent に

$$(3) \quad -\frac{d\psi}{dt} = \max_{f \in F} [\gamma(f) + Q(f)\psi(t)]$$

が成立する。(3)はいわゆる Bellman の最適性の原理であるが、Bellman は (3) を導いたわけではなく、(3) を以て最適条件と定義した。

更に、Miller は、homogeneous Markov (Q が t を explicit に含まない) の場合について、 π が最適であることと π が最大原理を満たすことが同値であることを証明し、また、piecewise constant な最適政策の存在定理を与えた。

Martin-Löf ([4]) は、finite state space, compact action space での non-homogeneous Markov の場合 ($Q(t, w)$), を扱い、 $Q(t, w)$, $\gamma(t, w)$ は共に有界かつリプシッツ条件を満たし、 $Q(t, w)$ は t に関して L -可測を仮定して、最適政策の存在を示した。更に詳細な結果として、periodic case では最適な periodic policy の存在を、また homogeneous case では最適な homogeneous policy の存在を示した。

Rykov ([8]) は、finite state, finite action での homogeneous Markov の場合、時間区間 $[0, +\infty)$ にわたっての平均 (単位時間当りの) 期待利得を Max ならしめる

問題を考之, (3) の analogy を導き, 本 stationary な最適政策の存在を示した.

2.2 ある非線形型確率微分方程式系

$$(4) \quad \frac{dX(\omega, t)}{dt} = f(t, X(\omega, t), u(\omega, t)) + \sigma(X(\omega, t), u(\omega, t)) \frac{dZ(\omega, t)}{dt}$$

$$(t \in \mathbb{R}^1, X \in \mathbb{R}^n, u \in \mathbb{R}^m)$$

Ω ; sample space

T ; 任意に fix

$X(\omega, t)$; Ω 上の確率過程

$\Sigma(t)$; $\{X(\cdot, \tau), 0 \leq \tau \leq t\}$ が可測となる Ω 上の最小の σ -field

$\mathcal{J}(T)$; $[0, T]$ 上の Borel field

$$\tilde{\Sigma}(T) \equiv \mathcal{B}\left(\bigcup_{0 \leq t \leq T} (\Sigma(t) \times \mathcal{J}(T))\right)$$

X^n ; \mathbb{R}^n のすべての compact subset からなる空間

つぎの仮定をおく.

- 1° 各 $\omega \in \Omega$ に対し, $u(\omega, \cdot)$ が $[0, T]$ 上で L -可測, 各 $t \in [0, T]$ に対し, $u(\cdot, t)$ が Ω 上で $\Sigma(t)$ 可測
 $u(\cdot, \cdot)$ は $\Omega \times [0, T]$ 上で $\tilde{\Sigma}(T)$ 可測
- 2° $A \equiv [0, T] \times C$ かつ C は \mathbb{R}^n の closed subset

3° control space $U(t, x)$ は $A \rightarrow X^m$ なる mapping で,
各 t に対し, x に関して連続 (Hausdorff distance で),
各 x に対し, t に関して L -可測

4° $f(t, x, u) = (f_1, f_2, \dots, f_n)$ は $A \times \bigcup_{(t,x) \in A} U(t, x) \rightarrow \mathbb{R}^n$ の mapping で,
各 t に対し, (x, u) に関して連続,
各 (x, u) に対し, t に関して L -可測

5° cost function $f_0(t, x, u) : A \times \bigcup_{(t,x) \in A} U(t, x) \rightarrow \mathbb{R}^1$
は, 各 t に対し, (x, u) に関して連続, 各 (x, u) に対し
 t に関して L -可測

6° $\tilde{f} = (f_0, f)$ とおくとき,

$$\|\tilde{f}(t, x, u)\| \leq m(t) \quad \text{for } (t, x, u) \in A \times \bigcup_{(t,x) \in A} U(t, x)$$

なる $[0, T]$ 上 L -可積分関数 $m(t)$ が存在する。

7° $Q(t, x) \equiv \left\{ (v^0, v) ; v^0 \geq f_0(t, x, u), \right.$
 $\left. v = f(t, x, u) + \sigma(x, u) \frac{dx}{dt}, u \in U(t, x) \right\}$

が殆んどすべての $\omega \in \Omega$ に対して 凸集合である。

$t_1(\omega), t_2(\omega)$ を $t_1(\omega) \leq t_2(\omega) \leq T$ w.p.1 なる random

time とし,

$$E \left[\int_{t_1(\omega)}^{t_2(\omega)} f_0(t, x(\omega, t), u(\omega, t)) dt \right] \rightarrow \text{Min}$$

ならしめる。

定理 2.1 (Furukawa)

$1^\circ \sim 7^\circ$ の仮定のもとで, いかなる non-empty な complete class の中でも最適は control が存在する.

参考文献

- [1] D. Blackwell ; Discounted dynamic programming,
Ann. Math. Statist., 36 (1965), 226-235.
- [2] N. Furukawa ; A Markov decision process with non-stationary transition laws, Bulletin of Math. Statist.,
13 (1968), 41-52.
- [3] R.A. Howard ; Dynamic programming and Markov processes, Technology Press of M. I. T and John Wiley,
(1960)
- [4] A. Martin-Löf ; Optimal control of a continuous-time Markov chain with periodic transition probabilities,
Operations Research, 15 (1967), 872-881
- [5] B.L. Miller ; Finite state continuous time Markov decision processes with a finite planning horizon,
SIAM Jour. on Control, 6 (1968), 266-280

- [6] J.L. Snell ; Application of martingale system theorems,
Trans. Amer. Math. Soc., 73 (1952), 293 - 312
- [7] R.E. Strauch ; Negative dynamic programming,
Ann. Math. Statist., 37 (1966), 871 - 890
- [8] V.V. Rykov. ; Markov decision processes with finite state
and decision spaces, Theory of Prob. and its Appl.,
11 (1966), 302 - 311