

Markov game with expected average reward criterion

新潟大・理 田 中 謙 輔

§1. 問題の定式化について

ここで非協力 N 人マルコフ・ゲームを $(2N+3)$ 個の組 $(S, A^1, A^2, \dots, A^N, \rho, f, r^1, r^2, \dots, r^N)$ で与えられる: (1) S は可分な距離空間でゲームの状態空間, (2) A^i は i player の行動空間と呼ばれるコンパクトな距離空間, (3) ρ は $S \times A$ の上で定義されている正值関数 $\rho(s, \bar{a})$, $\bar{a} = (a^1, a^2, \dots, a^N) \in \prod_{i=1}^N A_i = A$, (4) f は $s \in S$ と $\bar{a} \in A$ に対応する $(S, \beta(S))$ の上の確率測度 $f(\cdot | s, \bar{a})$, $\beta(S)$ は S の上の距離によって導入される Borel field である, (5) r^i は i player の利得率関数と呼ばれる $S \times A$ の上で定義されている実数値関数.

このゲームでは各 player が各時点でシステムの状態を観測し, 互に相談又は協力することなく現時点の状態 $s \in S$ のみで確率的に行動 $a^i \in A^i$ を選択する. この結果, i player

は利得 $r^i(s, \bar{a})$, $\bar{a} = (a^1, \dots, a^n) \in A$ を得る。そして ッステムは スルコフ過程に従って新しい状態に移って行く。このとき各 player の最適化問題は 各人の平均期待利得を最大にすることである。

ここでは各 player は現在の状態 $s \in S$ から $P(A^i)$ の中への可測写像である戦略のみを用いることとする、ただし $L P(A^i)$ は可測空間 $(A^i, \mathcal{B}(A^i))$ の上のすべての確率測度の集合とする。ゲーム ッステムの現時点までの歴史に関係していないうきょうな戦略は定常戦略と呼ばれている。初期状態 $s \in S$ と多重戦略 $\bar{\mu} = (\mu^1, \mu^2, \dots, \mu^n)$ に対して時刻 T までの player の期待利得は次のように与えられている。

$$\phi^i(s, T, \bar{\mu}) = E_{\bar{\mu}} \left[\int_0^T r^i(s_t, \bar{a}_t) dt \mid s_0 = s \right],$$

ただし s_t, \bar{a}_t は時刻 t における ッステムの状態と多重行動を示しており、 $E_{\bar{\mu}}$ は $\bar{\mu}$ によって決定される確率測度による平均を示している。さらに次のような記号を導入する

$$\phi^i(s, \bar{\mu}) = \overline{\lim_{T \rightarrow \infty}} \frac{\phi^i(s, T, \bar{\mu})}{T}$$

このとき、もし

$$\phi^i(s, \bar{\mu}) = \sup_{\alpha^i} \phi^i(s, (\mu^i, \alpha^i))$$

となる $\bar{\mu} = (\mu^1, \mu^2, \dots, \mu^n) \in \prod_{i=1}^n P(A^i) = P(A)$ が存在すれば、この $\bar{\mu}$ はマルコフ・ゲームの平衡点と呼ばれ、さらに各要素 μ^i は i player の平衡戦略と呼ばれている、ただし $(\mu^i, \alpha^i) = (\mu^1, \dots, \mu^{i-1}, \alpha^i, \mu^{i+1}, \dots, \mu^n)$.

2. 準備と記号

$C(A^i)$ を A^i 上の実数値連続関数の全体とする。 $F(A^i)$ を $(A^i, \beta(A^i))$ 上のすべての signed measure の集合とし、 $C(A^i)$ の双対空間から弱位相を導入する。このとき $F(A^i)$ は linear Hausdorff locally convex topological space となり $P(A^i)$ は $F(A^i)$ の compact convex metric subspace となる。 $\prod_{i=1}^n F(A^i)$ も linear Hausdorff locally convex topological space となり、 $P(A) = \prod_{i=1}^n P(A^i)$ は compact convex metric subspace となるので、 $P(A)$ は可分になっている。

さらに次のような仮定をとする

(A1)(i) $\phi(s, \bar{\alpha})$ は各 $\bar{\alpha} \in A$ に対して $\beta(S)$ 可測で、 $\forall s \in S$ に対して、 A 上で連続となっている。すなはち次のような正数 M が存在する：

$$0 < f(s, \bar{a}) \leq M \quad \forall s \in S, \forall \bar{a} \in A.$$

(ii) $f(\cdot | s, \bar{a})$ は、各 $s \in S \times \bar{a} \in A$ に対して、
 $(S, \mathcal{B}(S))$ 上の確率測度である。各 $\bar{a} \in A$ と $A \in \mathcal{B}(S)$ に
 対して、 f は $\mathcal{B}(S)$ 可測で、各 $A \in \mathcal{B}(S)$ と
 $s \in S$ に対して、 A 上で連続となつてゐる。さらには
 次のような条件をみたしてゐる： $\forall s \in S, \bar{a} \in A,$
 $A \in \mathcal{B}(S),$

$$0 \leq f(A | s, \bar{a}) \leq 1,$$

$$f(\{s\} | s, \bar{a}) = 0,$$

$$f(S | s, \bar{a}) = 1.$$

(AZ) $r^*(s, \bar{a})$ は、各 $\bar{a} \in A$ に対して、 $\mathcal{B}(S)$ 可測
 で、各 $s \in S$ に対して、 A 上の連続関数となつて
 いる。さらに次のような条件をみたす正数 N が存
 在する：

$$|r^*(s, \bar{a})| \leq N, \quad \forall s \in S, \bar{a} \in A.$$

ここで次のような記号も導入する： $P(A) \Rightarrow \bar{\mu}$ に対して、

$$\begin{aligned}
 Q(\Lambda, s, \bar{\mu}) &= \int_A \phi(s, \bar{a}) q(\Lambda | s, \bar{a}) d\bar{\mu}(\bar{a}) \\
 &= \int_{A^1} \cdots \int_{A^N} \phi(s, a^1, \dots, a^N) q(\Lambda | s, a^1, \dots, a^N) \prod_{i=1}^N d\mu^i(a^i), \\
 r^i(s, \bar{\mu}) &= \int_A r^i(s, \bar{a}) d\bar{\mu}(\bar{a}) \\
 &= \int_{A^1} \cdots \int_{A^N} r^i(s, a^1, \dots, a^N) \prod_{i=1}^N d\mu^i(a^i),
 \end{aligned}$$

上の等式は Fubini の定理より結論されている。このとき次の微分方程式は唯一 \rightarrow の解 $\phi^{\bar{\mu}}(t, \Lambda | s)$ a.e. $t \geq 0$ が存在する = λ が知られてる: a.e. $t \geq 0$, $\Lambda \in \mathcal{P}(S)$, $s \in S$, $\bar{\mu} \in \mathcal{P}(A)$

に対して、

$$\begin{aligned}
 \frac{d\phi^{\bar{\mu}}(t, \Lambda | s)}{dt} &= - \int_{\Lambda} \phi(x, \bar{\mu}) \phi^{\bar{\mu}}(t, dx | s) + \\
 &\quad + \int_S Q(\Lambda, x, \bar{\mu}) \phi^{\bar{\mu}}(t, dx | s), \tag{2.1}
 \end{aligned}$$

$$\phi^{\bar{\mu}}(0, \Lambda | s) = \delta(s, \Lambda) \quad \forall \bar{\mu} \in \mathcal{P}(A),$$

ただし

$$\delta(s, \Lambda) = \begin{cases} 1 & s \in \Lambda \\ 0 & s \notin \Lambda \end{cases}.$$

さらに、この解 $\phi^{\bar{\mu}}(t, \lambda | s)$ は次のような性質をみたしている：

- (1) $\phi^{\bar{\mu}}$ は各 $t \in [0, \infty)$ と $s \in S$ に対して, $(S, \beta(S))$ 上の確率測度である,
- (2) $\phi^{\bar{\mu}}$ は各 $s \in S$ と $\lambda \in \beta(S)$ に対して, t の絶対連続関数である,
- (3) $\phi^{\bar{\mu}}$ は各 $t \in [0, \infty)$ と $\lambda \in \beta(S)$ に対して $\beta(S)$ 可測になっている。

$M(S)$ は S 上の有界なボレル可測な実数値関数の全体からなる集合とすると, $M(S) \ni u$ に対して $\|u\| = \sup_s |u(s)|$ なる位相の導入によって $M(S)$ は Banach 空間となる。この "多重戦略" $\bar{\mu} \in P(A)$ に対して, その線形写像 $T_{\bar{\mu}}(\bar{\mu})$ と $A(\bar{\mu})$: $M(S) \rightarrow M(S)$ を次のように定義する: $u \in M(S)$ に対して

$$T_{\bar{\mu}}(\bar{\mu})u(s) = \int_S u(x) \phi^{\bar{\mu}}(t, dx | s)$$

$$A(\bar{\mu})u(s) = -\phi(s, \bar{\mu})u(s) + \int_S u(x) Q(dx, s, \bar{\mu})$$

このとき, $T_{\bar{\mu}}(\bar{\mu})$ と $A(\bar{\mu})$ は次の性質をみたしている:

$$(T1) \quad \lim_{t \rightarrow 0^+} T_t(\bar{\mu}) u(s) = T_0(\bar{\mu}) u(s) = u(s)$$

$$(T2) \quad T_t(\bar{\mu}) T_s(\bar{\mu}) u(s) = T_{t+s}(\bar{\mu}) u(s)$$

この性質をみたす $\{T_t(\bar{\mu}) : t \in [0, \infty)\}$ は 推移確率 $\hat{P}^{\bar{\mu}}$ に 対応する semi-group と呼ばれてる。

Lemma 1. 任意の多重戦略 $\bar{\mu} \in P(A)$ と $u \in M(S)$ は

対して

$$\frac{d T_t(\bar{\mu}) u(s)}{dt} = T_t(\bar{\mu}) A(\bar{\mu}) u(s)$$

が成立する。

Proof.

$$\begin{aligned} \frac{d T_t(\bar{\mu}) u(s)}{dt} &= \int_S u(x) \frac{d \hat{P}^{\bar{\mu}}(t, dx | s)}{dt} \\ &= \int_S u(x) (-\hat{P}(x, \bar{\mu})) \hat{P}^{\bar{\mu}}(t, dx | s) + \\ &\quad + \int_S u(x) \int_S Q(dx, dy, \bar{\mu}) \hat{P}^{\bar{\mu}}(t, dy | s) \\ &= \int_S \{-\hat{P}(x, \bar{\mu}) u(x) + \int_S u(y) Q(dy, x, \bar{\mu}) \hat{P}^{\bar{\mu}}(t, dx | s)\} \\ &= T_t(\bar{\mu}) A(\bar{\mu}) u(s). \end{aligned}$$

Lemma 2. 任意の多重戦略 $\bar{\mu} \in P(A)$ と $u \in M(S)$ は対して

$$u(s) = \int_0^\infty e^{-\alpha t} T_t(\bar{\mu}) (\alpha I - A(\bar{\mu})) u(s) dt$$

が成立する, ただし α は正の実数で, I は恒等写像である.

Proof.

$$\begin{aligned} & \int_0^\infty e^{-\alpha t} T_t(\bar{\mu}) (\alpha I - A(\bar{\mu})) u(s) dt \\ &= \int_0^\infty \left\{ \alpha e^{-\alpha t} T_t(\bar{\mu}) u(s) - e^{-\alpha t} \frac{d T_t(\bar{\mu}) u(s)}{dt} \right\} dt \\ &= \int_S u(x) \int_0^\infty \frac{d}{dt} \left(e^{-\alpha t} \hat{P}^{\bar{\mu}}(t, dx | s) \right) dt \\ &= u(s) \end{aligned}$$

上のような事実と Ky, Fan の不動点定理より次のよろづ補助定理が成立することが示される.

Lemma 3. $\alpha > 0$ に対して, 次の式をみたす多重戦略 $\bar{\mu}_* = (\mu_*^1, \mu_*^2, \dots, \mu_*^n) \in P(A)$ と $u(\mu_*^{\widehat{i}}) \in M(S)$ が存在する:

$$\begin{aligned} d u(\mu_*^{\widehat{i}}) &= \max_{\alpha^i} \{ r^i(s, (\mu_*^{\widehat{i}}, \alpha^i)) + A(\mu_*^{\widehat{i}}, \alpha^i) u(\mu_*^{\widehat{i}})(s) \} \\ &= r^i(s, \bar{\mu}_*) + A(\bar{\mu}_*) u(\mu_*^{\widehat{i}})(s) \end{aligned} \quad (2.2)$$

たゞ $\hat{\mu}_*^i = (\mu_*^1, \dots, \mu_*^{i-1}, \mu_*^{i+1}, \dots, \mu_*^N) \in \prod_{j \neq i} P(A^j)$.

§3. 割引因子をもたないマルコフ・ゲームにおける平衡点の存在について

この章では仮説 (A1), (A2) に加えて, さらには次のよろくな仮説を課すことにする.

(A3) 次のような条件をみたす正数 γ が存在する:

$$0 < \gamma \leq p(s, \bar{a}) \quad \forall s \in S, \forall \bar{a} \in A.$$

(A4) 次のような条件をみたすある状態 $s_0 \in S$ と正数 $\beta \in (0, 1)$ が存在する:

$$q(\{s_0\} | s, \bar{a}) > \beta > 0 \quad \forall \bar{a} \in A, s \neq s_0.$$

このとき (A3) の条件から (2.1) の解は honest 性をもつている, i.e.,

$$\bar{P}(t, S | s) = 1 \quad \forall s \in S, \bar{\mu} \in P(A).$$

今次のような新しい函数 $\bar{p}(s, \bar{a})$ と確率測度 $\bar{P}(\cdot | s, \bar{a})$ を作る:

$$\bar{P}(s, \bar{a}) = \begin{cases} p(s, \bar{a}) - \beta \gamma & s \neq s_0 \\ p(s, \bar{a}) & s = s_0 \end{cases}$$

任意の $\Lambda \in \beta(S)$ に対して

$$\bar{f}(\Lambda | s, \bar{a}) = \begin{cases} f(\Lambda | s, \bar{a}), & s = s_0 \\ \frac{p(s, \bar{a}) f(\Lambda | s, \bar{a})}{P(s, \bar{a})}, & s \neq s_0, s_0 \notin \Lambda \\ \frac{p(s, \bar{a}) f(\Lambda | s, \bar{a}) - \beta\gamma}{P(s, \bar{a})}, & s \neq s_0, s_0 \in \Lambda \end{cases}$$

上のような作りから $P(s, \bar{a}) \geq \bar{f}(\cdot | s, \bar{a})$ は仮説 (A1), (A2) を満たしてるので "Lemma 3 より" 次のような多重戦略 $\bar{\mu}_* = (\bar{\mu}_*^1, \bar{\mu}_*^2, \dots, \bar{\mu}_*^n) \in P(A)$ と $U(\bar{\mu}_*) \in M(S)$ が存在する: $\beta\gamma > 0$ に対して,

$$\begin{aligned} \beta\gamma U(\bar{\mu}_*)(s) &= \max_{\alpha_i} \{ r^i(s, (\bar{\mu}_*^i, \alpha^i)) + \bar{A}(\bar{\mu}_*, \alpha^i) U(\bar{\mu}_*)(s) \} \\ &= r^i(s, \bar{\mu}_*) + \bar{A}(\bar{\mu}_*) U(\bar{\mu}_*)(s), \end{aligned} \quad (3.1)$$

ただし

$$\bar{A}(\bar{\mu}) U(s) = -P(s, \bar{\mu}) U(s) + \int_S U(x) \bar{Q}(dx, s, \bar{\mu}),$$

$$\bar{Q}(\Lambda, s, \bar{\mu}) = \int_A P(s, \bar{a}) \bar{f}(\Lambda | s, \bar{a}) d\bar{\mu}(\bar{a}).$$

このとき (3.1) の中の \bar{A} を A に書き改めるにによって (3.1) 式は次のようになります: $\forall s \in S$ に対して,

$$\beta\gamma U(\bar{\mu}_*)(s_0) = \max_{\alpha_i} \{ r^i(s, (\bar{\mu}_*^i, \alpha^i)) + A(\bar{\mu}_*, \alpha^i) U(\bar{\mu}_*)(s) \} \quad (3.2)$$

同様にして、次式が得られる

$$\beta \delta u(\hat{\mu}_*^i)(s_0) = r^i(s, \bar{\mu}_*) + A(\bar{\mu}_*) u(\hat{\mu}_*^i)(s). \quad (3.3)$$

Theorem. 假設 (A1) ~ (A4) のもとで、割引因子をもたない非協力 N 人マルコフ・ゲームは平衡点をもつてゐる。

Proof. $g(\hat{\mu}_*^i) = \beta \delta u(\hat{\mu}_*^i)(s_0)$ とおくと、(3.2) は次のようすに書ける: $\forall \alpha^i \in P(A^i)$,

$$g(\hat{\mu}_*^i) \geq r^i(s, (\hat{\mu}_*^i, \alpha^i)) + A(\hat{\mu}_*^i, \alpha^i) u(\hat{\mu}_*^i)(s). \quad (3.4)$$

(3.4) の両辺に $T_*(\hat{\mu}_*^i, \alpha^i)$ を作用し、honest 性を用いると

$$g(\hat{\mu}_*^i) \geq T_*(\hat{\mu}_*^i, \alpha^i) r^i(s, (\hat{\mu}_*^i, \alpha^i)) + T_*(\hat{\mu}_*^i, \alpha^i) A(\hat{\mu}_*^i, \alpha^i) u(\hat{\mu}_*^i)(s)$$

を得る。

上式に Lemma 1 を用いて、次に 0 から T まで積分すると
とから次式に達する:

$$T g(\hat{\mu}_*^i) \geq \int_0^T T_*(\hat{\mu}_*^i, \alpha^i) r^i(s, (\hat{\mu}_*^i, \alpha^i)) dt + T_*(\hat{\mu}_*^i, \alpha^i) u(\hat{\mu}_*^i)(s) - u(\hat{\mu}_*^i)(s).$$

両辺を T で割って、次に $T \rightarrow \infty$ として上極限をとると

$$\begin{aligned} g(\hat{\mu}_*^i) &\geq \overline{\lim}_{T \rightarrow \infty} \frac{\phi^i(s, T, (\hat{\mu}_*^i, \alpha^i))}{T} \\ &= \phi^i(s, (\hat{\mu}_*^i, \alpha^i)) \end{aligned} \quad (3.5)$$

が得られる。

(3.3) から同様な議論によつて

$$g(\bar{\mu}_*^i) = \phi^i(s, \bar{\mu}_*) \quad (3.6)$$

が得られる。

(3.5) と (3.6) から

$$\phi^i(s, \bar{\mu}_*) = \max_{\alpha^i} \phi^i(s, (\bar{\mu}_*^i, \alpha^i))$$

が成立し、 $\bar{\mu}_* \in P(A)$ が マルコフ・ケームの平衡点になっている。よつて証明は完成された。

References

1. Beneš, V.E., Existence of optimal strategies based on specific information for a class of stochastic decision problems, SIAM J. Control and Optimization, 8(1970).
2. Fan,Ky, Fixed points and minimax theorems in locally convex topological linear spaces, Proc. Nat. Acad. Sci., USA, 38(1952).
3. Feller, W., On the integro-differential equations of purely discontinuous Markov processes, Trans. Amer. Math. Soc., 48(1940).
4. Lai, H.C. and Tanaka,K., Non-cooperative n-person game with a stopped set, to appear in J. Math. Analy. Appl.
5. Tanaka,K. and Lai, H.C., Two-person zero-sum game with a stopped set, to appear in J. Math. Analy. Appl.
6. Tanaka,K. and Homma,H., Continuous time non-cooperative n-person Markov games, Bull. Math. Statist., 15(1978).
7. Tanaka,K. and Wakuta,K., On continuous time Markov games with countable state space, J. O.R. Japan, 21(1978).
8. Tanaka,K. and Homma,H., On the learning algorithm of 2-person zero-sum Markov game, Bull. Math. Statist., 19 (1980).
9. Tanaka,K., Markov Game, Proceedings of symposium on functional analysis and applications, at National Tsing Hua University, Taiwan, Republic of China, 1980.
10. Tanaka,K., Stochastic differential games, Report on Research under Financial support of the National Sciences, Council, Republic of China, 1980.