

ある sequential assignment problemについて

大阪府立大学総合科学部 中井 遼

1. Introduction

S を state space, P を stochastic transition function とする stationary Markov chain を (S_m, P) とするとき、この Markov chain の上での sequential stochastic assignment problem $\{A_n, S, P, r, \{f_{t+1}, s, p\}, \beta\}$ を考える。但し A_n は action space であり、二つの disjoint set $\{a_1, \dots, a_m\}$ と $\{p\}$ の和とし、あらわせる。また全ての state $s \in S$ に対して non-negative random variable X_s の値 x_s を観測することが出来る。この確率変数は独立で、分布関数 $F(x)$ は既知とする。 $r(a, x)$ は観測値 x に対する action a を用いたときの reward とし $E[r(a, x)] < \infty$ であると仮定する。 β は discount factor であり $0 \leq \beta \leq 1$ とする。

現在まで考えられている sequential stochastic assignment problem と同様に、次の形を sequential assignment problem を考える。([1][3][4][5][7])。まず position (A_n, s) を現在の state が s であり、そのときの action space が A_n である状態とする。position が (A_n, s) であるとき、確率変数 X_s の

実現値を観測して action $a \in A_n$ 中から選ぶ。たとえば $\{a_1, \dots, a_m\}$ のうちの a_i ならば、immediate reward $r(a_i, x_s)$ を得る。このとき次の stage t の position は (A_n, t) である。 $A_n = \{a_1, \dots, a_{n-1}, a_n, \dots, a_m\}$ たり A は Markov chain (S, P) の next stage t である。また a が action "p" に等しい時は、immediate reward は 0 であり、次の stage t の position は (A_n, t) である。即ち action "p" は pass することと同じである。その回数には制限はない。以上の様な場合の、total expected discounted reward を最大にする問題を考える。すなはち action "a_i" を取るとは別ゆゑ sequential assignment problem で、action a_i を観測値 x と assign すること等しい。

2. Problem formulation.

Sequential stochastic assignment problem $\{A_n, S, P, r\}$ が与えられ、現在の position が (A_n, s) であるとする。optimal strategy の下での value $E V(A_n, s)$ 、また、観測値 x を知り、 E という条件付の。この problem の value $E V(A_n, s|x)$ を求めよう。よし β と γ とする。dynamic programming formulation によれば次の recursive equation が得られる。

$$V(A_n, s) = E_x V(A_n, s|x)$$

(1)

$$V(A_n, s|x) = \max_{1 \leq i \leq n} \max_{a \in A_n} \{r(a_i, x) + \beta \int V(A_n^t, t) P(s, dt)\},$$

$$\beta \int U(A_{n+1}) P(x, dx) \gamma$$

但し、ここでおいて $A_n = \{a_1, \dots, a_m\}^{\text{per}}$, $A_n' = \{a_1, \dots, a_{l-1}, a_l\}$, $a_m \notin A_n'$ とする。(1)式と同様によく知られた Markov decision process と同様に optimal strategy の存在と value の存在が示される。ここで $\gamma \propto T$, optimal strategy & value が求められる。 reward function に条件を付ける。即ち $r(a, x) = r(a) b(x)$ とあらわせる場合を考える。また state が x であるときの確率変数を X_0 の代りに $T_2(X_0)$ を考えることにより $r(a, x) = r(a)x$ と考える。また A_n に対し $r(a_1) \geq \dots \geq r(a_m)$ を仮定する。このとき optimal strategy 及びその value について次の事が成り立つ。

Proposition 1. 全ての $a \in S$ に対して次の value function の列 $\{g_i(b)\}_{i=1, \dots, m}$, $\{h_i(b)\}_{i=1, 2, \dots}$ が存在する。即ち

$$g_1(b) \geq g_2(b) \geq g_3(b) \geq \dots \geq g_n(b) \geq \dots \geq 0$$

$$h_1(b) \geq h_2(b) \geq h_3(b) \geq \dots \geq h_n(b) \geq \dots \geq 0$$

である。position (A_n, s) において次の事が成り立つ。

1) 觀測値の値が x であるときの optimal strategy は

$g_i(b) \leq x < g_{i-1}(b)$ のとき i -th action a_i を選ぶ ($i=1, \dots, n$)
 $0 \leq x < g_n(b)$ のとき action "p" を選ぶ事である。

2) position (A_n, s) のときの value は

$$U(A_n, s) = \sum_{i=1}^m r(a_i) h_i(b)$$

3) $\{g_i(b)\}_{i=1,2,\dots}$ 及び $\{h_i(b)\}_{i=1,2,\dots}$ は次の recursive equation によって求められることが出来る。

$$g_i(b) = \beta \int h_i(t) P(b, dt).$$

$$h_i(b) = \int_{g_{i-1}(b)}^{\infty} g_{i-1}(x) dF_b(x) + \int_0^{g_i(b)} x dF_b(x) + \int_0^{g_i(b)} g_i(x) dF_b(x),$$

$$\text{又 } h_i(b) = \int_{g_i(b)}^{\infty} x dF_b(x) + \int_0^{g_i(b)} g_i(x) dF_b(x), \quad g_i(b) = \beta \int h_i(t) P(b, dt)$$

$$g_0(b) = \infty.$$

さて、次に $P_{t_0} \in P(T_0 \leq t) = P(s+t)$ であるより確率変数とする。今 $s+t$ へ収束するとき X_{t_0} 及び T_{t_0} がそれぞれ X_0 及び T_0 へ確率収束すると仮定する。この仮定の下で次の事が成り立つ。

Corollary 1 $g_i(b)$ ($b \in S, i=1, 2, \dots$) は S 上で連続である。

以下では上の事を仮定して話を進める。

3. Special case ($S \subset R^n \times R'$)

この節で stochastic transition function P が、現在の state s のみならず、観測値 x の値を depend する場合を考える。即ち前節の definition と従がえは、 $S \subset R^n \times R'$ で、 $S \ni (s, x)$ を考える。また、 $dF_{(s, x)}(y) = I_x$ (I_x は indicator function) かつ $P((s, x), (t, y)) = P(s, x, t) F_t(y)$ である場合を想定する。この問題では、Proposition 1 における函数 $g_i(s, x)$ が、よ

1) 簡単を試すが出来る。但し、ここでは次の様な
仮定を設ける。即ち $P(s, x, t)$ は x に関して stochastically
increasing かつ $\partial P(s, x, t) / \partial x^2 \geq 0$, 更に $F_0(x)$ は x に関して
stochastically increasing であるとする。

今 $U(A_n, s) = E v(A_n, s, X_s)$ かつ $U(A_n, s | x) = U(A_n, s, x)$
と定義すれば、二点の値は、optimal strategy の下の value
及び、条件付の value であり、(1)式を用いて次の recursive
equation を得ることが出来る。

$$U(A_n, s) = \int_0^\infty U(A_n, s | x) d F_0(x)$$

$$(2) \quad U(A_n, s | x) = \max \left\{ \max_{1 \leq i \leq n} \left\{ v(a_i) x + \beta \int_0^x U(A_{n-i}, t) d P(s, dt) \right\}, \right. \\ \left. \beta \int_x^\infty U(A_{n-i}) d P(s, dt) \right\},$$

但し $A_n = \{a_1, \dots, a_n\} \cup \{p\}$ かつ $A_n^i = \{a_1, \dots, a_{i-1}, a_{i+1}, \dots, a_n\} \cup \{p\}$

このとき次の事が成り立つ。

Proposition 2. 全ての $s \in S$ に対して 次の形を満たす

$$t_1(s) \geq t_2(s) \geq \dots \geq t_n(s) \geq \dots \geq 0.$$

且つ $t_i(s)$ が存在する。

このとき position (A_n, s) における optimal strategy β の value
は次の形にあるわせる。

1) 觀測値が x のとき optimal strategy は

$f_i(s) \leq x < f_{i+1}(s)$ ($i=1, \dots, n$) のとき i -th action a_i を選ぶ

$0 \leq x < f_n(s)$ のとき action "p" を選ぶ事である。

2) value $U(A_{n,s})$ は

$$U(A_{n,s}) = \sum_{i=1}^n r(a_i) t_i(\omega) \text{ である}。$$

3) $\{t_i(\omega)\}_{i=1,2,\dots}$ は次の recursive equation を満足する。

$$t_n(\omega) = \int_{d_{n+1,0}}^{\infty} t_{n+1}(\omega, x) dF_0(x) + \int_{d_{n,s}}^{d_{n+1,s}} x dF_0(x) + \int_0^{d_{n,s}} t_n(\omega, x) d\bar{F}_0(x)$$

$$t_n(\omega, x) = \beta \int t_n(t) P(\omega, x, dt),$$

但し $t_0(\omega) = \infty$, $-\infty \cdot 0 = 0 \cdot \infty = 0 \Rightarrow d_{n,s}$ は $t_{n+1}(x) = x$ の unique root である。

ここで、前節の Proposition 1 で述べた $\{g_i(\omega, x)\}_{i=1,2,\dots}$ 及び $\{h_i(\omega, x)\}_{i=1,2,\dots}$ と Proposition 2 における $\{f_i(\omega)\}_{i=1,2,\dots}$ 及び $\{t_i(\omega, x)\}_{i=1,2,\dots}$ の関係は、次のようにあるわす事が出来る。

$$t_n(\omega) = \int h_n(\omega, x) dF_0(x)$$

$$\begin{aligned} t_n(\omega, x) &= \beta \int t_n(t) P(\omega, x, dt) = \beta \int \int h_n(t, y) dF_t(y) P(\omega, x, dt) \\ &= g_n(\omega, x) \end{aligned}$$

上の Proposition 2 を用いて次の二つの Corollary を得る。

Corollary 2. $t_i(\omega)$ は ω に関して増加関数である。

Corollary 3. $t_i(\omega, x)$ は x に関して増加関数でありまた上に凸な関数である。

4. Variants and example

4.1 2節で考えた問題は N -stage problem とて考えれば、
Proposition 1 と 同様の結果が得られ、 recursive equation と π^* を
計算して optimal strategy & value を得ることが出来た。

Corollary 4. 全ての state $s \in S$ に対し、 固定列 $\{g_{i,N}(s)\}_{i=1 \dots N}$
及 $U\{h_{i,N}(s)\}_{i=1 \dots N}$ で π^* の状況下での存在する。

$$g_{1,N}(s) \geq g_{2,N}(s) \geq \dots \geq g_{N-1,N}(s) \geq g_{N,N}(s) = 0.$$

$$h_{1,N}(s) \geq h_{2,N}(s) \geq \dots \geq h_{N-1,N}(s) \geq h_{N,N}(s) = 0.$$

position at $N(A_n, s)$ であると次の事が成り立つ。(但し
 $N(A_n, s)$ は、残り stage の数を N と、 position at (A_n, s) である
時を表す。またそのときの value は $U_N(A_n, s)$ である。)

1) 繼測値の値が x であるときの optimal strategy は

$$g_{i,N}(s) \leq x < g_{i+1,N}(s) \text{ のとき } i\text{-th action } a_i \in \pi^* \quad (i=1, \dots, N)$$

$0 \leq x < g_{N,N}(s)$ のとき action "p" をとることである。

2) value $U_N(A_n, s)$ は

$$U_N(A_n, s) = \sum_{i=1}^n r(a_i) h_{i,N}(s) \text{ である。}$$

3) $\{g_{i,N}(s)\}_{i=1 \dots N}$ 及 $U\{h_{i,N}(s)\}_{i=1 \dots N}$ は π^* の recursive
equation を満足する。

$$g_{1,N}(s) = \beta \int h_{1,N}(x) P(s, dx)$$

$$h_{i,N}(s) = \int_{g_{i-1,N-1}(s)}^{\infty} g_{i,N-1}(x) d\bar{F}_0(x) + \int_{g_{i,N-1}(s)}^{g_{i,N}(s)} x d\bar{F}_0(x) + \int_0^{g_{i,N}(s)} g_{i,N}(x) d\bar{F}_0(x)$$

$$\text{但し } h_{1,N}(w) = \int_{g_{1,N}(w)}^{\infty} x dF_0(x) + \int_0^{f_{1,N}(w)} f_{1,N}(t) dF_0(x)$$

$$g_{1,N}(w) = \int h_{1,N}(t) P(a, dt). \text{ である。}$$

4.2. $W \in \text{space } \mathcal{Q}$ の値 ϵ と 3 parameter, X を実数値確率変数とし, $f(x|w)$ を $w=w$ であるときの X の conditional g . p.d.f. とする。 $g(w)$ を W の prior g.p.d.f. とし, $g(w|x)$ を W の posterior g.p.d.f. とする。

現在まで m 個の値 $x_1 = x_1, \dots, x_m = x_m$ を observe したときの posterior distribution を $G_m(w|x_1, \dots, x_m)$ として次の仮定を設ける。(但し $G_0(w) = G(w)$ とする。) 1) $\mu_w = E[X|w] < \infty$ ($w \in \mathcal{Q}$) かつ $\mu = E[\mu_w] < \infty$. 2) 与えられた x_1, \dots, x_m に対する w に対する充分統計量 \bar{x}_m が存在し. 3) \bar{x}_m は \bar{x}_{m-1} 及び x_m によって recursively 生成される。 4) $E[\bar{w}|\bar{x}_m]$ は \bar{x}_m の non-decreasing function であり 5) $P(X_m \geq z | \bar{x}_{m-1})$ は \bar{x}_{m-1} に関する stochastically increasing とする。

このとおり 3 節と同様に、次の様な sequential stochastic assignment problem を考える。 $A_n = \{a_1, \dots, a_n\} \neq \emptyset$, $S = N \times \mathbb{R}^d$
 $\Rightarrow (m, \bar{x}_m)$, $r(a, x) = r(a) \cdot x$ ($r(a_1) \geq \dots \geq r(a_m)$), $F_{(m, \bar{x}_m)}$
 $= F(x|\bar{x}_m) \Leftrightarrow P((m, \bar{x}_m), x, (m', y)) = I_{(m+1, \bar{x}_{m+1}(\bar{x}_m, x))}$,
 但し $(\bar{x}_{m+1}(\bar{x}_m, x))$ は \bar{x}_m 及び w , $(m+1)$ stage における観測値 x によって生成される充分統計量とする。このとおり Proposition 2.

\rightarrow より optimal strategy 及び value が得られる。この結果は Nakai [3] における解と一致する。

4.3 Random termination である場合の sequential stochastic assignment problem を考えよ。この場合の problem は $\{A_n, S, P, r, \beta\}$ である。特に $S = S_1 \cup \{S^*\}, P(T=S^*) = 1, P(X_{S^*}=0) = 1$, 但し T は $P(T \leq t) = P(S^* \leq t)$ である確率変数である。また S^* は absorbing state であり, S^* の観測値の値は常に 0 である。この場合, absorbing state S^* に入れば, reward は常に 0 であり, そのときこの問題は終わると考えられる。

4.4 この sequential stochastic assignment problem は game 的な扱いも考えられる。即ち, 二人の player が, 連続出現する観測値から m 個を選択。それに対して, action space A_n の m 個の action を取らせる場合を考える。但し player I は観測値の選択と同時に, どの action をとるかという選択を含むせて行ない。player II は観測値の選択のみを行なう。但しこのとき, 両 player が, ここで考えた problem と異なり, pass する回数は制限されているものとする。また player I は maximizing player として player II は minimizing player として行動をとるものとする。このゲームについては Nakai [5] の解が述べられていて、同様に optimal strategy

及 V value を得 $3 = 2 \times 2^2 - 3.$

References

- [1] C. Derman, G.J. Lieberman and S.M. Ross "A sequential stochastic assignment problem" Management Science vol 18, p349 - 355, 1972.
- [2] G.H. Hardy, J.E. Littlewood and G. Polya "Inequalities" Cambridge University Press, 1934.
- [3] T. Nakai "Optimal assignment for a random sequence with an unknown parameter" Journal of Information & Optimization Science, vol 1, p214~228, 1980.
- [4] T. Nakai "Sequential stochastic assignment problem with rejection" Journal of Information & Optimization Science, vol 2, p169~180, 1981.
- [5] T. Nakai "A time sequential game related to the sequential stochastic assignment problem", Journal of Operations Research Society of Japan, vol 25, No. 2, 1982
- [6] M. Sakaguchi "A sequential assignment problem for randomly arriving jobs", Rep. Stat. Appl. Res. JUSE, vol 19, p99~109, 1972.