

折扣模型最优策略的结构

董泽清 刘 克

(中国科学院应用数学研究所)

摘要

本文研究了折扣马尔可夫决策规划(以下简记为MDP)最优策略的结构。证明了：任给一策略 $\pi^* = (\pi_0^*, \pi_1^*, \dots, \pi_n^*, \pi_{n+1}^*, \dots)$ ，若它是 β 折扣最优的，则随机平稳策略 $\pi_0^{*\infty} = (\pi_0^*, \pi_0^*, \dots)$ 也是 β 折扣最优的；对任何 $n(\geq 1)$ ，我们也给出了随机平稳策略 $\pi_n^{*\infty} = (\pi_n^*, \pi_n^*, \dots)$ 也是 β 折扣最优的充分条件。还证明了：任给一随机平稳策略 $\pi_0^{\infty} = (\pi_0, \pi_0, \dots)$ ，若它是 β 折扣最优的，则必可分解为若干个 β 折扣最优的(决定性)平稳策略，而且这些平稳策略的凸组合就得到原来的 π_0^{∞} 。因此，较完满地解决了折扣模型最优策略的结构。

§ 1 引 言

我们沿用[1]的概念和记号，不重述。本文所研究的MDP： $\{S, (A(i), i \in S), b, r, V_\beta\}$ ，其中状态空间 S 、每个状态可用的行动(决策)集 $A(i)$ 均为可列集，状态转移律簇 q 是时齐的，报酬函数 r 是有界的(界记作 M)，以及折扣目标

$$V_\beta(\pi; i) = \sum_{t=0}^{\infty} \beta^t \sum_{\substack{a \in A(j) \\ j \in S}} P_\pi \{Y_t = j, A_t = a \mid Y_0 = i\}, \quad \pi \in \Pi, i \in S, \quad (1)$$

其中 $\beta \in (0, 1)$ 是给定的。

一个策略 $\pi^* \in \Pi$ ，若对任给的 $i \in S$ 、 $\pi \in \Pi$ 均有

$$V_\beta(\pi^*; i) \geq V_\beta(\pi; i) \quad (2)$$

则称 π^* 是 β 折扣最优的或 π^* 是 β 折扣最优策略，简称为最优的或最优策略。

若一个(决定性)马氏策略 $\pi^* = (f_0^*, f_1^*, f_2^*, \dots)$ 是最优的，则对同一折扣因子 β ， $f_0^{*\infty}$ 也是最优的，这乃一熟知的事实。设一个策略 $\pi^* = (\pi_0^*, \pi_1^*, \pi_2^*, \dots, \pi_n^*, \pi_{n+1}^*, \dots)$ 是最优的，那么对任何非负整数 n ，时刻 n 的(随机)决策规则 π_n^* 构成的平稳策略 $\pi_n^{*\infty}$ 对同一 β 是否也是最优的？除 $n = 0$ 之外，回答是否定的；我们也给出了 $\pi_n^{*\infty}$ 也是最优策略的条件；以及给出了修改 $\pi_n^{*\infty}$ 的办法，使修改后的 π_n^* 所构成的平稳策略，对同一 β 也是最优的。

Chitgopekar[3]已证明：若 f^∞ 与 g^∞ 是两个(决定性) β 折扣最优策略，则它们的随机化(即凸组合)构成的随机平稳策略也是 β 折扣最优的。反问题：若随机平稳策略 $\pi_0^\infty = (\pi_0^\infty, \pi_0^\infty, \pi_0^\infty, \dots)$ 是 β 折扣最优的，令

$$A'(i) = \{a; \pi_0(a \mid i) > 0\}, \quad i \in S, \quad (3)$$

则任取一个 $f(i) \in A'(i)$ ， $i \in S$ ，问这样定义的 f 所构成的平稳策略 f^∞ 对同一 β 是否也是最

*1984年2月20日收到。

优的？回答是肯定的。即最优随机平稳策略，可分解若干个决定性平稳策略，它们对同一 β 都是最优的。

这样，我们就完满地解决了折扣模型最优策略的结构。粗略地说，一个策略是最优的，当而且仅当，在每个时刻，对每个可能到达的状态都必须选用最优的行动。

§ 2. 最优策略的性质

令 B 为定义在 S 上的可列维列向量空间，对任给的 $V, V' \in B$ ，若对每个 $i \in S$ ，均有 $V(i) \geq V'(i)$ ，则记作 $V \geq V'$ ；若 $V \geq V'$ 且存在 $i \in S$ 使 $V(i) > V'(i)$ ，则记作 $V > V'$ 。令

$$B = \{V : -(1-\beta)^{-1}M\mathbf{e} \leq V \leq (1-\beta)^{-1}M\mathbf{e}, V \in B\}, \quad (4)$$

其中 M 是报酬函数 r 的绝对值之界， $e \in B$ ，且每个分量均为1的列向量。 F 为全体决策函数所成之集。

对任给的 $f \in F, V \in B$ ，定义映射 T_f 为

$$T_f V = r(f) + \beta Q(f)V, \quad (5)$$

其中 $r(f)$ 为一列向量，其第 i 个分量为 $r(i, f(i))$ ， $Q(f)$ 为一（次）随机矩阵，其第 (i, j) 个元素为 $q(j|i, f(i))$ ， $i, j \in S, \beta \in (0, 1)$ 。在下面的讨论中，我们总认为 β 是事先给定的。

我们不加证明地引用如下结果（证明可见〔1〕的第5章），把它们归并成一个引理。

引理1 对任给的 $f \in F, V \in B$ ，我们有：

- 1) T_f 把 B 映射到 B ；
- 2) T_f 为一单调映射，即对任给的 $V, V' \in B$ ，若 $V \geq V'$ ，则 $T_f V \geq T_f V'$ ；
- 3) T_f 有唯一不动点 $V_\beta(f^\infty) \in B$ ，即

$$T_f V_\beta(f^\infty) = V_\beta(f^\infty),$$

其中 $V_\beta(f^\infty)$ 的第 i 个分量 $V_\beta(f^\infty; i)$ 由(1)式给出，只要取 $\pi = f^\infty$ ：

- 1) 若 $T_f V \geq V$ ，则 $V_\beta(f^\infty) \geq V, V \in B$ ；
- 2) 设 $\pi' = (g_0, g_1, g_2, \dots) \in \prod_m^d, \pi \in \prod$ ，则

$$\begin{aligned} T_{\pi_0} T_{\pi_1} \cdots T_{\pi_n} V_\beta(\pi) &= V_\beta(g_0, g_1, \dots, g_n, \pi) \\ &\rightarrow V_\beta(\pi'), \text{ 当 } n \rightarrow \infty. \end{aligned}$$

定理1 设 $\tau^* = (\tau_0^*, \pi_0^*, \tau_1^*, \pi_1^*, \dots) \in \Pi$ 是 β 折扣最优的，且 $\tau_0^* = f \in F$ ，则 f^∞ 对同一 β 也是最优的。

证 对任给的 $i \in S$ ，我们均有

$$V_\beta(\tau^*; i) = r(i, f(i)) + \beta \sum_{j \in S} q(j|i, f(i)) V_\beta(\tau^*(i, f(i)); j).$$

其中

$$\tau^*(i, f(i)) = (\pi'_0, \pi'_1, \pi'_2, \dots) \in \Pi$$

具有

$$\pi'_t(\cdot | s_0, a_0, s_1, a_1, \dots, s_t) \equiv \tau_{t+1}^*(\cdot | i, f(i), s_0, a_0, s_1, a_1, \dots, s_t),$$

$$a_n \in A(s_n), s_n \in S, n = 0, 1, 2, \dots, t,$$

即 $\tau^*(i, f(i))$ 是由 π 往后推移一个时段所构成的策略。而去掉的那个时段，系统必处于状态 i 。

选用行动 $f(i)$ ，由于 π^* 是最优的，则

$$\begin{aligned} V_\beta(\pi^*; i) &\leq r(i, f(i)) + \beta \sum_{j \in S} q(j | i, f(i)) V_\beta(\pi^*; j) \\ &= [T_f V_\beta(\pi^*)](i), \quad i \in S \end{aligned}$$

即

$$V_\beta(\pi^*) \leq T_f V_\beta(\pi^*),$$

由引理 (1.4) 有 $V_\beta(\pi^*) \leq V_\beta(f^\infty)$ ，由 π^* 是最优的，故 $V_\beta(\pi^*) = V_\beta(f^\infty)$ 。 ■

系 (Blackwell [2]) 的定理 2) 若 $\pi = (f_0, f_1, f_2, \dots) \in \Pi_m^d$ 是最优的，则 f_n^∞ 对同一 β 也是最优的。

设 $\pi = (f_0, f_1, \dots, f_n, \pi_{n+1}, \pi_{n+2}, \dots)$ 是 β 折扣最优的，那么， f_n^∞ 对同一 β 是否也是最优的？回答是否定的。反例如下：

例 取 $S = \{0, 1\}$, $A(0) = \{1, 2\}$, $A(1) = \{1\}$, $q(1 | 0, 1) = q(1 | 0, 2) = q(1 | 1, 1) = 1$, $r(0, 1) = 1$, $r(0, 2) = r(1, 1) = 0$ 。

显然 $F = \{f, g\}$, 其中 $f(i) = 1$, $i = 0, 1$; $g(0) = 2$, $g(1) = 1$ 。

由于 F 是有限集，故存在平稳策略是最优的（如见 [1] 的第 2 章）。易见 f^∞ 是最优的（对任给的 $\beta \in (0, 1)$ 均成立），其最优报酬向量为

$$V_\beta(f^\infty; 0) = \sum_{t=0}^{\infty} \beta^t \sum_{j \in S} P_{\pi^*}(Y_t = j | Y_0 = 0) r(j, f(j)) = 1,$$

$$V_\beta(f^\infty; 1) = \sum_{t=0}^{\infty} \beta^t \sum_{j \in S} P_{\pi^*}(Y_t = j | Y_0 = 1) r(j, f(j)) = 0;$$

以及 $V_\beta(g^\infty; i) = 0$, $i \in S$ ，即 g^∞ 不是最优的（对任给的 $\beta \in (0, 1)$ 均成立）。若取 $\pi = (f, g, g, \dots)$ ，我们有

$$V_\beta(\pi; 1) = 0$$

$$V_\beta(\pi; 0) = r(0, f(0)) + \beta \sum_{j \in S} g(j | 0, f(0)) V_\beta(g^\infty) = 1,$$

故

$$V_\beta(\pi) \equiv V_\beta(f^\infty),$$

即 π 是最优的，但 g^∞ 不是最优的。因此，一般来讲， $\pi = (f_0, f_1, \dots, f_n, \pi_{n+1}, \pi_{n+2}, \dots)$ 是最优的，对同一 β , f_n^∞ 不一定是最优的。下面我们给出一个充分条件，它保证 f_n^∞ 仍是最优的。为此，先引入在一个策略下“可实现的历史”这一概念。

定义 设 $\pi = (\pi_0, \pi_1, \pi_2, \dots) \in \Pi$, 对任给的 $t (\geq 1)$, 若存在历史 $h_{t-1}(\pi) \equiv (i_0, a_0, i_1, a_1, \dots, i_{t-1}, a_{t-1})$ 使

$$\pi_0(a_0 | i_0) q(i_1 | i_0, a_0) \pi_1(a_1 | h_0(\pi), i_1) q(i_2 | i_1, a_1) \dots \cdot$$

$$q(i_{t-1} | i_{t-2}, a_{t-2}) \pi_{t-1}(a_{t-1} | h_{t-2}(\pi), i_{t-1}) > 0$$

则称 $h_{t-1}(\pi)$ 为使用策略 π ，于 $t = 0$ 时从状态 i_0 出发的条件下，直到 $t - 1$ 时刻可实现的历史，其中 $a_n \in A(i_n)$, $i_n \in S$, $n = 0, 1, \dots, t - 1$ 。对给定的 π, t ，这种可实现的历史之全体

$\{h_{t-1}(\pi)\}$ 记做 $H_{t-1}(\pi)$ 。对给定的 $j \in S$, $\pi \in \Pi$ ，若存在某个可实现的历史 $h_{t-1}(\pi) \equiv (i_0, a_0, i_1, a_1, \dots, i_{t-1}, a_{t-1}) \in H_{t-1}(\pi)$ 使 $g(j | i_{t-1}, a_{t-1}) > 0$ ，则称使用策略 π ，于 $t = 0$ 时从 i_0

出发的条件下, 通过可实现的历史 $h_{t-1}(\pi)$ 在时刻 t 可达状态 j , 简称为用策略 π 、 $t=0$ 从 i_0 出发, 于时刻 t 可达 j .

定理 2 若 $\pi = (f_0, f_1, \dots, f_n, \pi_{n+1}, \pi_{n+2}, \dots) \in \Pi$ 是 β 折扣最优的, 且对任给的 $j \in S$ 必存在 $i_0 \in S$ 使得, 在使用策略 π 、 $t=0$ 从状态 i_0 出发, 于时刻 n 可达状态 j , 则 f_n^∞ 对同一 β 也是最优的.

证 对任给的 $r (r=1, 2, \dots, n)$, $i_k \in S (k=0, 1, \dots, r-1)$, 令

$$\begin{aligned} & r\pi(i_0, f_0(i_0), i_1, f_1(i_1), \dots, i_{r-1}, f_{r-1}(i_{r-1})) \equiv r_\pi(h_{r-1}(\pi)) \\ & = (f_r, f_{r+1}, \dots, f_n, r\pi'_{n+2}, \dots) \in \Pi, \end{aligned} \quad (6)$$

其中

$$\begin{aligned} & r\pi'_{n+r} (\cdot | s_0, f_r(s_0), \dots, s_{n-r}, f_n(s_{n-r}), s_{n-r+1}, a_{n-r+1}, \dots, s_{n-r+t}) \\ & = \pi_{n+r} (\cdot | h_{r-1}(\pi), s_0, f_r(s_0), \dots, s_{n-r}, f_n(s_{n-r}), s_{n-r+1}, a_{n-r+1}, \dots, s_{n-r+t}) \\ & a_k \in A(s_k), k=n-r+1, n-r+2, \dots, n-r+t, \\ & s_j \in S, j=0, 1, \dots, n-r+t, t=r+1, r+2, \dots, \end{aligned}$$

具有 $h_{r-1}(\pi) = (i_0, f_0(i_0), i_1, f_1(i_1), \dots, i_{r-1}, f_{r-1}(i_{r-1}))$.

由于 π 是 β 折扣最优的, 及定理 1 知

$$V_\beta(\pi) \equiv V_\beta(f_0^\infty) \geq V_\beta(\pi'), \pi' \in \Pi. \quad (7)$$

由 (6)、(1) 及 (7), 对任给的 $h_{n-1}(\pi) = (i_0, f_0(i_0), \dots, i_{n-1}, f_{n-1}(i_{n-1}))$ 我们有

$$\begin{aligned} & V_\beta(\pi; i_0) = r(i_0, f_0(i_0)) + \beta \sum_{i_1} q(i_1 | i_0, f_0(i_0)) V_\beta(\pi(i_0, f_0(i_0)); i_1) \\ & = r(i_0, f_0(i_0)) + \beta \sum_{i_1} q(i_1 | i_0, f_0(i_0)) r(i_1, f_1(i_1)) \\ & + \beta^2 \sum_{i_2} q(i_2 | i_0, f_0(i_0)) \sum_{i_1} q(i_1 | i_0, f_0(i_0)) V_\beta(\pi(i_0, f_0(i_0)); i_1) \\ & = \dots = r(i_0, f_0(i_0)) + \beta \sum_{i_1} q(i_1 | i_0, f_0(i_0)) r(i_1, f_1(i_1)) \\ & + \beta^2 \sum_{i_2} q(i_2 | i_0, f_0(i_0)) \sum_{i_1} q(i_1 | i_0, f_0(i_0)) r(i_2, f_2(i_2)) + \dots \\ & + \beta^{n-1} \sum_{i_1} q(i_1 | i_0, f_0(i_0)) \sum_{i_2} q(i_2 | i_1, f_1(i_1)) \dots \sum_{i_{n-1}} q(i_{n-1} | i_{n-2}, f_{n-2}(i_{n-2})) \\ & \cdot r(i_{n-1}, f_{n-1}(i_{n-1})) \\ & + \beta^n \sum_{i_1} q(i_1 | i_0, f_0(i_0)) \dots \sum_{i_n} q(i_n | i_{n-1}, f_{n-1}(i_{n-1})) V_\beta(\pi(i_0, f_0(i_0)); i_n) \\ & \leq r(i_0, f_0(i_0)) + \beta \sum_{i_1} q(i_1 | i_0, f_0(i_0)) r(i_1, f_1(i_1)) \\ & + \beta^2 \sum_{i_2} q(i_2 | i_0, f_0(i_0)) \sum_{i_1} q(i_1 | i_0, f_0(i_0)) r(i_2, f_2(i_2)) + \dots \\ & + \beta^{n-1} \sum_{i_1} q(i_1 | i_0, f_0(i_0)) \dots \sum_{i_{n-1}} q(i_{n-1} | i_{n-2}, f_{n-2}(i_{n-2})) r(i_{n-1}, f_{n-1}(i_{n-1})) \\ & + \beta^n \sum_{i_1} q(i_1 | i_0, f_0(i_0)) \dots \sum_{i_n} q(i_n | i_{n-1}, f_{n-1}(i_{n-1})) V_\beta(\pi; i_n), \end{aligned} \quad (8)$$

上式中未标明的求和范围, 均是在 S 上求和 (下同), 若用向量、矩阵来写, 再用引理 1.5 有

$$V_\beta(\pi) \leq V_\beta(f_0, f_1, \dots, f_{n-1}, \pi) \leq V_\beta(\pi).$$

因此,(8)式全成立等式,并得到

$$\begin{aligned} & \beta^n \sum_{i_1} q(i_1 | i_0, f_0(i_0)) \cdots \sum_{i_n} q(i_n | i_{n-1}, f_{n-1}(i_{n-1})) V_\beta(\pi^{(h_{n-1}(\pi))}; i_n) \\ &= \beta^n \sum_{i_1} q(i_1 | i_0, f_0(i_0)) \cdots \sum_{i_n} q(i_n | i_{n-1}, f_{n-1}(i_{n-1})) V_\beta(\pi; i_n), \end{aligned}$$

由于 π 是最优的,因此

$$V_\beta(\pi^{(h_{n-1}(\pi))}; i_n) \leq V_\beta(\pi; i_n), \quad i_n \in S.$$

故,对任一可实现历史 $h_{n-1}(\pi)$,于时刻 n 可达的那些 i_n 必有

$$V_\beta(\pi^{(h_{n-1}(\pi))}; i_n) = V_\beta(\pi; i_n). \quad (9)$$

由所设,对每个 $j \in S$,必存在 $i_0 \in S$,使得在用策略 π , $t=0$ 从状态 i_0 出发的条件下,于时刻 n 可达 j .因此,对每个 i_n ,均必存在可实现历史 $h_{n-1}(\pi)$,使

$$V_\beta(\pi^{(h_{n-1}(\pi))}; i_n) = V_\beta(\pi; i_n).$$

对每个 i_n ,用这样的历史,我们就确定了一个 β 折扣最优策略

$$(f_n, \pi'_{n+1}, \pi'_{n+2}, \dots) \equiv \pi' \quad \Pi, \quad (10)$$

再由定理1,就得到 f_n^∞ 对同一 β 也是最优的.■

从(9)式我们看到,若不要求对所有状态出发同时达最优,而 $\pi = (f_0, f_1, \dots, f_n, \pi_{n+1}, \pi_{n+2}, \dots)$ 是最优的,凡是使用策略 π ,从某个初始($t=0$ 时)状态出发,于时刻 n 可达的那些状态,然后在 $t=0$ 改为从这些可达状态出发,用策略 π' (由(10)给出),对同一 β 它也是最优的.但对于使用策略 π ,从任何初始状态出发,于时刻 n 均不可达的那些状态,在 $t=0$ 时,改为从这些不可达状态出发、用策略 π' ,就无法保证它还是最优的.这实质上就是Bellman的“最优化原理”对我们所研究的模型成立.

如何修改 f_n ,使修改后的 f_n 能构成一个平稳最优策略?我们有下结果.

定理3 若 $\pi = (f_0, f_1, \dots, f_n, \pi_{n+1}, \pi_{n+2}, \dots)$ 是 β 折扣最优的,令

$$A^*(i) = \{a; r(i, a) + \beta \sum_j q(j | i, a) V_\beta(\pi; j) = V_\beta(\pi; i), \quad a \in A(i)\}.$$

$$\overline{B} = \left\{ i; \begin{array}{l} \text{存在 } i_0 \in S, \text{ 使用策略 } \pi, t=0 \text{ 从 } i_0 \text{ 出发, 于时刻 } n \\ \text{可达状态 } i, \quad i \in S \end{array} \right\}$$

且定义 f_n^* 为

$$f_n^*(i) = \begin{cases} f_n(i), & \text{当 } i \in \overline{B}, \\ a, & \text{当 } i \notin \overline{B}, \quad a \in A^*(i) \end{cases} \quad i \in S$$

则 f_n^{∞} 对同一 β 也是最优的.

这个定理的证明,从下面的定理7立即可得.为了简单起见,当 $i \in \overline{B}$ 时可取

$$f_n^*(i) = f_0(i).$$

一个问题:定理1中之 π^* 的初如决策规则若是随机的 π_0^* ,会有什么结果?我们先叙述一引理.

对任给的 $i \in S$, $\pi_0(\cdot | i)$ 是 $A(i)$ 上的概率分布.类似于(5)我们定义一个映射 \mathcal{F}_{π_0} 为:

$$\begin{aligned} [\mathcal{F}_{\pi_0} V](i) &= \sum_{a \in A(i)} \pi_0(a | i) [r(i, a) + \beta \sum_j q(j | i, a) V(j)], \\ V &\in B, \quad i \in S. \end{aligned} \quad (11)$$

若用向量、矩阵符号来写,(11) 可写作

$$\mathcal{J}_{\pi_0} V = \pi_0 r + \beta \pi_0 q V, \quad V \in \mathbf{B}, \quad (12)$$

其中 $\mathcal{J}_{\pi_0} V$, $\pi_0 r$, $\beta \pi_0 q V$ 均为列向量, 其第 i 个分量分别为 $[\mathcal{J}_{\pi_0} V](i)$, $\sum_{a \in A(i)} \pi_0(a|i) r(i, a)$, 以及 $\beta \sum_{a \in A(i)} \pi_0(a|i) \sum_j q(j|i, a) V(j)$.

引理 2 对任给的随机规则 $\pi_0, V \in \mathbf{B}$, 按 (11)(或 (12)) 定义的 \mathcal{J}_{π_0} , 我们有

- 1) \mathcal{J}_{π_0} 是把 \mathbf{B} 映射到 \mathbf{B} 的单调映射;
- 2) 若 $\mathcal{J}_{\pi_0} V > V$, 则 $V_\beta(\pi_0^\infty) > V$;
- 3) 若 $\mathcal{J}_{\pi_0} V \geq V$, 则 $V_\beta(\pi_0^\infty) \geq V$;
- 4) 若 $\mathcal{J}_{\pi_0} V < V$, 则 $V_\beta(\pi_0^\infty) < V$;
- 5) 若 $\mathcal{J}_{\pi_0} V \leq V$, 则 $V_\beta(\pi_0^\infty) \leq V$.

证 1) $\forall i \in S$, 由 (11)

$$[\mathcal{J}_{\pi_0} V](i) \leq \sum_{a \in A(i)} \pi_0(a|i) [|r(i, a)| + \beta \sum_j q(j|i, a) |V(j)|] \\ \leq M + \beta M(1 - \beta)^{-1} = M(1 - \beta)^{-1},$$

故 \mathcal{J}_{π_0} 把 \mathbf{B} 映射到 \mathbf{B} . 设 $V, V' \in \mathbf{B}$ 且 $V \geq V'$, 由 (12),

$$\mathcal{J}_{\pi_0} V - \mathcal{J}_{\pi_0} V' = \pi_0 r + \beta \pi_0 q V - (\pi_0 r + \beta \pi_0 q V') = \beta \pi_0 q (V - V') \geq 0,$$

即 \mathcal{J}_{π_0} 是单调映射.

2) 由 1), $\mathcal{J}_{\pi_0} \mathcal{J}_{\pi_0} V \equiv \mathcal{J}_{\pi_0}^2 V \geq \mathcal{J}_{\pi_0} V > V$, 即 $\pi_0 r + \beta \pi_0 q \pi_0 r + \beta^2 \pi_0 q V > V$, 用归纳法, 对任何正整数 n 有 $\mathcal{J}_{\pi_0}^n V \geq \mathcal{J}_{\pi_0} V > V$, 即

$$\pi_0 r + \beta \pi_0 q \pi_0 r + \beta^2 \pi_0 q \pi_0 q \pi_0 r + \dots + \beta^{n-1} \pi_0 q \pi_0 q \dots \pi_0 q \pi_0 r + \\ \underbrace{\beta^n \pi_0 q \pi_0 q \dots \pi_0 q \pi_0 q}_n V \geq \pi_0 r + \beta \pi_0 q V > V, \quad (13)$$

最左端最后一项的绝对值 $\leq \beta^n (1 - \beta)^{-1} M \cdot e \rightarrow 0$, 当 $n \rightarrow \infty$ 时 (按分量分别取极限), 其中 e 是每个元素均为 1 的列向量. 而

$$V_\beta(\pi_0^\infty) = \pi_0 r + \beta \pi_0 q \pi_0 r + \beta^2 \pi_0 q \pi_0 q \pi_0 r + \dots + \beta^{n-1} \pi_0 q \pi_0 q \dots \pi_0 q \pi_0 r + \\ \underbrace{\beta^n \pi_0 q \pi_0 q \dots \pi_0 q \pi_0 q}_n \pi_0 r + \dots$$

上式右端前 n 项与 (13) 式最左端的前 n 项完全相同, 在 (13) 式两端令 $n \rightarrow \infty$, 得

$$V_\beta(\pi_0^\infty) \geq \mathcal{J}_{\pi_0} V > V.$$

3)、4) 及 5) 的证明完全类似于 2) 的证明

定理 4 设策略 $\pi^* = (\pi_0^*, \pi_1^*, \pi_2^*, \dots) \in \Pi$ 是 β 折扣最优的, 则 $\pi_0^* = (\pi_0^*, \pi_1^*, \pi_2^*, \dots)$ 对同一 β 也是最优的.

证 由于 $V_\beta(\pi^*) = \pi_0^* r + \beta \pi_0^* q \pi_1^* r + \beta^2 \pi_0^* q \pi_1^* q \pi_2^* r + \dots + \beta^n \pi_0^* q \pi_1^* q \dots \pi_{n-1}^* q \pi_n^* r + \dots$
 $= \pi_0^* r + \beta \pi_0^* q [\pi_1^* r + \beta \pi_1^* q \pi_2^* r + \dots + \beta^{n-1} \pi_1^* q \pi_2^* q \dots \pi_{n-1}^* q \pi_n^* r + \dots]$
 $= \pi_0^* r + \beta \pi_0^* q V_\beta(\pi^*(\cdot, \cdot, \dots)),$

按分量来写, $\forall i \in S$

$$V_\beta(\pi^*; i) = \sum_{a \in A(i)} \pi_0^*(a|i) r(i, a) + \sum_{a \in A(i)} \pi_0^*(a|i) \sum_j q(j|i, a) V_\beta(\pi^*(\cdot, \cdot, \dots); j),$$

其中 $\pi^*(\cdot, \cdot, \dots) = (\pi_0^*, \pi_1^*, \pi_2^*, \dots)$ 具有 $\pi_t^*(\cdot, s_0, a_0, s_1, a_1, \dots, s_t) = \pi_{t-1}^*(\cdot, \cdot, \dots, a_k, s_k, a_k, \dots, s_t)$, $s_i, a_i, \dots, s_t \in S$, $a_k \in A(s_k)$, $s_k \in S$, $k = 0, 1, \dots, t$, $t = 0, 1, 2, \dots$,

显然 $\pi^*(\cdot, \cdot, \dots) \in \Pi$ (但依赖于初始状态与初始行动), 由于 π^* 是最优的, 因此

$$V_\beta(\pi^*) \leq \pi_0^* r + \beta \pi_0^* q V_\beta(\pi^*) = \mathcal{J}_{\pi_0} V_\beta(\pi^*)$$

由引理2.3有 $V_\beta(\pi_0^*, \infty) \geq V_\beta(\pi^*)$, 由于 π^* 的最优性, 必有 $V_\beta(\pi_0^*, \infty) = V_\beta(\pi^*)$, 即 π_0^*, ∞ 对同一 β 也是最优的.

下面的定理5与6的证明, 分别类似于定理2与3的证明, 现仅叙述如下.

定理5 若 $\pi = (\pi_0, \pi_1, \dots, \pi_n, \pi_{n+1}, \dots) \in \Pi$ 是 β 折扣最优的, 且对任给的 $j \in S$, 必存在 $i_0 (i_0 \in S)$, 使得在用策略 π 时, $t=0$ 从 i_0 出发, 于时刻 n 可达状态 j , 则 π_n^∞ 对同一 β 也是最优的, 其中 π_n , 对每个 $j \in S$, 只要任取一通过 $h_{n-1}(\pi)$ 于时刻 n 可达 j 的历史, $\pi_n(\cdot | h_{n-1}(\pi), j)$ 就完全确定了.

定理6 若 $\pi = (\pi_0, \pi_1, \dots, \pi_n, \pi_{n+1}, \dots) \in \Pi$ 是 β 折扣最优的, 任给 $i \in S$, 若 $i \in \bar{B}$, 则任取一个元 $(\cdot | i)$ 为 $A^*(i)$ 上的概率分布, 则定义 π_n^* 为

$$\pi_n^*(\cdot | i) = \begin{cases} \pi_n(\cdot | h_{n-1}(\pi), i), & \text{当 } i \in \bar{B}, \\ \pi(\cdot | i), & \text{当 } i \notin \bar{B}, \end{cases}$$

则 π_n^*, ∞ 对同一 β 也是最优的, 其中 $\bar{B}, A^*(i)$ 的定义同定理3, $h_{n-1}(\pi)$ 为使用策略 π 的一个可实现历史、并通过它于时刻 n 可达 i .

上面诸定理的证明本质上未用到报酬函数 r 的有界性.

§ 3 最优随机平稳策略的分解

现在来回 § 1 中提出的第二个问题. 令

$$V_\beta^*(i) = \sup_{\pi \in \Pi} V_\beta(\pi; i), \quad i \in S,$$

即 V_β^* 为最优报酬向量, 我们有 (见 [1] 的第五章)

$$V_\beta^*(i) = \sup_{a \in A(i)} [r(i, a) + \beta \sum_j q(j | i, a) V_\beta^*(j)], \quad i \in S, \quad (14)$$

$A^*(i)$ 由定理3中所给出, 注意到那里的 π 是最优的, 所以 $V_\beta(\pi) \equiv V_\beta^*$. $A^*(i)$ 实为状态 i 可用的最优行动集 (若 $A^*(i)$ 非空时).

引理3 若存在 β 折扣最优策略 $\pi^* = (\pi_0^*, \pi_1^*, \dots)$, 则对一切 $i \in S$, $A^*(i)$ 均是非空的.

证 由于 π^* 是最优的, 对一切 $i \in S$, 由 (14) 式有

$$\begin{aligned} V_\beta^*(i) &= V_\beta(\pi^*; i) = \sum_{a \in A(i)} \pi_0^*(a | i) [r(i, a) + \beta \sum_j q(j | i, a) V_\beta(\pi^{*(i, a)}; j)] \\ &\leq \sum_{a \in A(i)} \pi_0^*(a | i) [r(i, a) + \beta \sum_j q(j | i, a) V_\beta^*(j)] \\ &\leq \sup_{a \in A(i)} [r(i, a) + \beta \sum_j q(j | i, a) V_\beta^*(j)] = V_\beta^*(i), \end{aligned} \quad (15)$$

因此, 上式全成立等式, 其中

$$V_\beta(\pi^{*(i, a)}; j) = \sum_{t=1}^{\infty} \beta^{t-1} \sum_{\substack{b \in A(k) \\ k \in S}} P_\pi \cdot Y_t = k, \quad A_t = b | Y_0 = i, A_0 = a, Y_1 = j \} r(k, b),$$

由 (14), 我们总有

$$r(i, a) + \beta \sum_j q(j | i, a) V_\beta^*(j) \leq V_\beta^*(i), \quad a \in A(i), \quad i \in S,$$

由 (15) 成立等式, 因此, 对一切使 $\pi_0^*(a_0 | i) > 0$ 的 a_0 均应有

$$r(i, a_0) + \beta \sum_j q(j | i, a_0) V_\beta^*(j) = V_\beta^*(i),$$

故, $A^*(i)$ 非空, $i \in S$.

定理 7 一个策略 $\pi = (\pi_0, \pi_1, \dots, \pi_n, \pi_{n+1}, \dots)$ 是 β 折扣最优的充要条件是对每个 $n \geq 0$, 时刻 n 的决策规则 π_n 局限在最优行动集上也构成一个概率分布, 即在使用策略 π 下, 通过任一可实现历史 $h_{n-1}(\pi)$, 于时刻 n 可达 j 的条件下, 选取行动 $a \in A^*(j)$ 的概率为 $\pi_n(a | h_{n-1}(\pi), j)$, 且有 $\sum_{a \in A^*(j)} \pi_n(a | h_{n-1}(\pi), j) = 1$.

证 必要性: 设 $\pi^* = (\pi_0^*, \pi_1^*, \dots, \pi_n^*, \dots)$ 是最优的, 由引理 3, 知 π_0^* 对任何 $i \in S$, $\pi_0^*(\cdot | i)$ 均为 $A^*(i)$ 上的概率分布, $A^*(i)$ 由定理 3 给出, 而且 (15) 全成立等式, 即

$$\begin{aligned} V_\beta^*(i) &= V_\beta(\pi^*; i) = \sum_{a \in A^*(i)} \pi_0^*(a | i) [r(i, a) + \beta \sum_j q(j | i, a) V_\beta(\pi^{*(i, a)}; j)] \\ &= \sum_{a \in A^*(i)} \pi_0^*(a | i) [r(i, a) + \beta \sum_j q(j | i, a) V_\beta^*(j)] \\ &= E_{\pi^*}[r(Y_0, A_0) | Y_0 = i] + \sum_{a \in A^*(i)} \pi_0^*(a | i) \beta \sum_j q(j | i, a) V_\beta(\pi^{*(i, a)}; j) \\ &= E_{\pi^*}[r(Y_0, A_0) | Y_0 = i] + \sum_{a \in A^*(i)} \pi_0^*(a | i) \beta \sum_j q(j | i, a) V_\beta^*(j). \end{aligned}$$

由于总有

$$V_\beta(\pi^{*(i, a)}; j) \leq V_\beta^*(j), \quad j \in S,$$

因此对一切使 $P_{\pi^*}\{\mathcal{A}_0 = a, Y_1 = j | Y_0 = i\} = \pi_0^*(a | i) q(j | i, a) > 0$ 的 (i, a, j) 均有

$$V_\beta(\pi^{*(i, a)}; j) = V_\beta^*(j), \quad j \in S, \quad a \in A^*(i), \quad i \in S.$$

同理可证明

$$\begin{aligned} V_\beta^*(i) &= V_\beta(\pi^*; i) = E_{\pi^*}[r(Y_0, A_0) | Y_0 = i] + \beta E_{\pi^*}[r(Y_1, A_1) | Y_0 = i] \\ &\quad + \beta^2 \sum_{a \in A^*(i)} \pi_0^*(a | i) \sum_j q(j | i, a) \sum_{b \in A^*(j)} \pi_1^*(b | i, a, j) \sum_k q(k | j, b) V_\beta(\pi^{*(i, a, j, b)}; k) \\ &= E_{\pi^*}[r(Y_0, A_0) | Y_0 = i] + \beta E_{\pi^*}[r(Y_1, A_1) | Y_0 = i] \\ &\quad + \beta^2 \sum_{a \in A^*(i)} \pi_0^*(a | i) \sum_j q(j | i, a) \sum_{b \in A^*(j)} \pi_1^*(b | i, a, j) \sum_k q(k | j, b) V_\beta^*(k), \end{aligned}$$

且当 $P_{\pi^*}\{\mathcal{A}_0 = a, Y_1 = j | Y_0 = i\} > 0$ 时, $\pi_1^*(\cdot | i, a, j)$ 仅是 $A^*(j)$ 上的概率分布, 以及当

$$P_{\pi^*}\{\mathcal{A}_0 = a, Y_1 = j, A_1 = b, Y_2 = k | Y_0 = i\} > 0$$

时有

$$V_\beta(\pi^{*(i, a, j, b)}; k) = V_\beta^*(k), \quad a \in A^*(i), \quad b \in A^*(j), \quad i, j, k \in S,$$

其中

$$V_\beta(\pi^{*(i, a, j, b)}; k) = \sum_{t=2}^{\infty} \beta^{t-2} E_{\pi^*}[r(Y_t, A_t) | h_1(\pi), Y_2 = k],$$

而 $h_1(\pi) = (i, a, j, b)$.

用归纳法我们可以证明: 当 $P_{\pi^*}\{\mathcal{A}_0 = a_0, Y_1 = i_1, A_1 = a_1, \dots, Y_n = i_n | Y_0 = i\} > 0$ 时, $\pi_n^*(\cdot | h_{n-1}(\pi^*), i_n)$ 仅是 $A^*(i_n)$ 上的概率分布, 当

$$P_{\pi^*}\{\mathcal{A}_0 = a_0, Y_1 = i_1, A_1 = a_1, \dots, Y_n = i_n, A_n = a_n, Y_{n+1} = i_{n+1} | Y_0 = i\} > 0$$

时有

$$V_\beta(\pi^{*(h_n(\pi^*))}; i_{n+1}) = V_\beta^*(i_{n+1}), \quad i_{n+1} \in S,$$

其中 $h_n(\pi^*) = (i, a_0, i_1, a_1, \dots, i_n, a_n)$, $a_t \in A^*(i_t)$, $i_t \in S$, $t = 1, 2, \dots, n$.

这就证明了, π^* 在各个时刻选取行动的规则局限在最优行动集上, 必为一概率分布.

充分性：设 $\pi^* = (\pi_0^*, \pi_1^*, \pi_2^*, \dots, \pi_n^*, \pi_{n+1}^*, \dots)$ 中的 $\pi_n^* (n \geq 0)$ 均是局限于最优行动集上的概率分布，它也为一个策略，对任给的 $i \in S$ ，

$$\begin{aligned} V_\beta(\pi^*; i) &= E_{\pi^*} \left[\sum_{t=0}^{\infty} \beta^t r(Y_t, A_t) \mid Y_0 = i \right] \\ &= E_{\pi^*} \left\{ \left[\sum_{t=0}^{N-1} \beta^t r(Y_t, A_t) + \beta^N V_\beta(\pi^*(Y_0, A_0, \dots, Y_{N-1}, A_{N-1}); Y_N) \right] \mid Y_0 = i \right\} \end{aligned}$$

其中 $\pi^*(Y_0, A_0, \dots, Y_{N-1}, A_{N-1})$ 的第 $n+1$ 个决策规则为

$$\begin{aligned} \pi_n^*(i, A_0, \dots, Y_{N-1}, A_{N-1}) &= \pi_{N+n}^*(\cdot | i, A_0, \dots, Y_{N-1}, A_{N-1}, Y_N, A_N, \dots, Y_{N+n}), \\ n &\geq 0, N \geq 1. \end{aligned}$$

而

$$\begin{aligned} &E_{\pi^*} \left\{ \left[\sum_{t=0}^{N-1} \beta^t r(Y_t, A_t) + \beta^N V_\beta^*(Y_N) \right] \mid Y_0 = i \right\} \\ &\leq E_{\pi^*} \left\{ \left[\sum_{t=0}^{N-1} \beta^t r(Y_t, A_t) + \beta^N (V_\beta(\pi^*(Y_0, A_0, \dots, Y_{N-1}, A_{N-1}); Y_N) + \frac{2M}{1-\beta}) \right] \mid Y_0 = i \right\} \\ &= E_{\pi^*} \left\{ \left[\sum_{t=0}^{N-1} \beta^t r(Y_t, A_t) + \beta^N V_\beta(\pi^*(Y_0, A_0, \dots, Y_{N-1}, A_{N-1}); Y_N) \right] Y_0 = i \right\} \\ &\quad + \beta^N \frac{2M}{1-\beta} = V_\beta(\pi^*; i) + \beta^N \frac{2M}{1-\beta}, \end{aligned}$$

即

$$V_\beta(\pi^*; i) + \beta^N \frac{2M}{1-\beta} \geq E_{\pi^*} \left\{ \left[\sum_{t=0}^{N-1} \beta^t r(Y_t, A_t) + \beta^N V_\beta^*(Y_N) \right] \mid Y_0 = i \right\}, \quad i \in S, N \geq 1, \quad (16)$$

由 $A^*(j)$ 的定义，我们有

$$\begin{aligned} &E_{\pi^*} \left\{ \left[\beta^{N-1} r(Y_{N-1}, A_{N-1}) + \beta^N V_\beta^*(Y_N) \right] \mid Y_0 = i \right\} \\ &= \beta^{N-1} \sum_{\substack{a \in A^*(j) \\ j \in S}} E_{\pi^*} \left\{ [r(j, a) + \beta V_\beta^*(Y_N)] \mid Y_0 = i, Y_{N-1} = j, A_{N-1} = a \right\} \\ &\quad \cdot P_{\pi^*} \{Y_{N-1} = j, A_{N-1} = a \mid Y_0 = i\} \\ &= \beta^{N-1} \sum_{\substack{a \in A^*(j) \\ j \in S}} \{r(j, a) + \beta \sum_k q(k | j, a) V_\beta^*(k)\} P_{\pi^*} \{Y_{N-1} = j, A_{N-1} = a \mid Y_0 = i\} \\ &= \beta^{N-1} \sum_{a, j} V_\beta^*(j) P_{\pi^*} \{Y_{N-1} = j, A_{N-1} = a \mid Y_0 = i\} = \beta^{N-1} E_{\pi^*} [V_\beta^*(Y_{N-1}) \mid Y_0 = i], \end{aligned}$$

将上式代入 (16)，则有

$$\begin{aligned} V_\beta(\pi^*; i) + \beta^N \frac{2M}{1-\beta} &\geq E_{\pi^*} \left\{ \left[\sum_{t=0}^{N-2} \beta^t r(Y_t, A_t) + \beta^{N-1} V_\beta^*(Y_{N-1}) \right] \mid Y_0 = i \right\} \\ &= \dots = E_{\pi^*} \{[r(Y_0, A_0) + \beta V_\beta^*(Y_1)] \mid Y_0 = i\} = V_\beta^*(i), \quad i \in S, N \geq 1, \end{aligned}$$

然后，令 $N \rightarrow \infty$ 得 $V_\beta(\pi^*; i) \geq V_\beta^*(i)$, $i \in S$ ，

由 $V_\beta^*(i)$ 是最优报酬 ($t=0$, 从 i 出发), 知 π^* 是最优的.

这个定理说明：一个策略是最优的，当且仅当在每个时刻的决策规则，对可达的那些状态都必须选用最优行动。因此，若

$$\pi_0^\infty = (\pi_0, \pi_1, \dots) \in \Pi_S,$$

是最优的随机平稳策略，则 $\pi_0(a|i) > 0$ 的 a ，必属于 $A^*(i)$, $i \in S$ ，即由(3)式定义的 $A'(i) \subset A^*(i)$, $i \in S$ 。因此，对每个 $i \in S$ ，任取一 $f(i) \in A'(i)$ ，则 f^∞ 对同一 β 也必须是最优的。即 π_0^* 必是若干个（可能是无限个）决策函数的凸组合而成，而且这些决策函数之每一个，都单独构成一个平稳策略是最优的（对同一 β ）。

Π^* , Π_m^{d*} , Π_m^* , Π_s^{d*} , Π_s^* 分别表示 Π , Π_m^d , Π_m , Π_s^d , Π_s 中的最优策略集。 $C(A)$ 表示集 A 的凸包，则我们已证明下结果。

定理 8 若MDP: $\{S, (A(i), i \in S), q, r, V_\beta\}$ 存在最优策略，必有 $C(\Pi_s^{d*}) = \Pi_s^* \subset \Pi^*$ ；类似地有 $C(\Pi_m^{d*}) = \Pi_m^* \subset \Pi^*$ ，这几 $C(\Pi_m^{d*})$ 是诸最优（决定性）马氏策略，各时刻决策规则的凸组合，所成的凸包。

从而我们较完满地搞清了折扣MDP最优策略的结构。

参 考 文 献

- [1] 董泽清，“马尔可夫决策规划”讲义，科学院应用数学所油印，1981年。
- [2] Blackwell, D., Discrete dynamic programming, Ann. Math. Statist. Vol. 33(1962), 719—726.
- [3] Chitgopekar, S.S., Denumerable state Markovian sequential control processes: On randomizations of optimal policies, Naval Res. Logist. Quart. Vol. 22(1975), 567—573.
- [4] 郭世贞，折扣目标最优策略中的最小方差策略问题，昆明工学院油印，1983年。