

# 可控馬尔可夫鏈的一种最优决策

吳滄浦

## 摘 要

本文研究了一种最优馬尔可夫控制系统,这种控制系统以統計規律依赖于决定序列的馬尔可夫鏈描述。我們称决定序列为决策。存在一具有下述性质的目标状态,一旦系統到达此状态,状态就不再改变。我們的目的是要选取一决策,使所有从每一初始状态出发最終到达此目标状态的概率都达到最大。我們先提出在平稳决策集合中求最优决策的决策迭代法。然后証明,此决策在包含平稳及不平稳决策的决策集合上也是最优的。

## 一、問題的提法

某些最优控制过程可以归結为一类最优可控馬尔可夫鏈的模型。Bellman 首先研究了这类模型的渐近性质<sup>[1,2]</sup>。Howard 研究了两种这类模型的最优决策的求解問題<sup>[3]</sup>。在一种模型中,他假定任何决策下的馬尔可夫鏈都是正则的,指标是依赖于极限轉移概率的某种平均期望效益。在另一种模型中,未来阶段的效益要“打折扣”,这时,指标是依赖于轉移概率的总期望效益。Eaton 和 Zadeh 研究了一种具有吸收型目标状态的最优可控馬尔可夫鏈<sup>[4]</sup>,指标是依赖于轉移概率的总期望代价。在他們的模型中,决策集合只包含平稳决策。Blackwell 和 Derman 还进一步研究了决策集合含有不平稳决策及随机决策时最优决策的性质<sup>[5,6,7,8]</sup>。

以最小化总期望代价为指标,可以描写实际問題中某些平均值的指标,如平均过渡历程時間、平均誤差等,但不能描写某些概率指标。例如,在 Понтрягин 所研究的一类最优追踪問題中<sup>[9]</sup>,被追点在相空間中的运动是一个馬尔可夫过程,追点的运动是可控的,最优控制的指标是追点命中被追点的概率。在这問題中,如果把点在相空間中的运动用格子点的离散形式表示,則追踪过程可归結为具有吸收型目标状态的可控馬尔可夫鏈。这时,最优指标不能用上述工作中的期望值指标表示,而应该以到达目标状态的概率(下称命中概率)作为指标。

本文将以最大化命中概率作为指标,研究在包含不平稳决策的决策集合中,具有吸收型目标状态的可控馬尔可夫鏈的最优决策。

具有吸收型目标状态的可控馬尔可夫鏈可描述如下。在時間的进程中,系統只在給定的离散时刻  $t_0, t_1, t_2, \dots$  发生状态的改变。系統所有可能的状态数是有限的,我們以  $i = 1, 2, \dots, n, n+1$  表示之。在状态改变后,系統的状态变成某一新状态的概率,只取决于改变前它所处的状态以及所取的控制。在任一时刻  $t_k (k = 0, 1, 2, \dots)$  之前,如系統处于状态  $i (i = 1, 2, \dots, n)$ ,則可供选择的控制  $d_i$  有  $l_i$  个,  $d_i$  选定后,相应的各轉移概率  $p_{ij}(d_i) (j = 1, 2, \dots, n, n+1)$  也就确定。  $i = n+1$  表示目标状态。从目标状态出发的轉移概率  $p_{n+1,i}$  則是固定的:  $p_{n+1,1} = p_{n+1,2} = \dots = p_{n+1,n} = 0, p_{n+1,n+1} = 1$ , 即

目標狀態是一個吸收狀態。

由此可見，在時刻  $t_k$ ，鏈的轉移概率矩陣  $P$  取決於所選取的一組控制  $d_1^k, d_2^k, \dots, d_n^k$ 。我們以向量  $\mathbf{d}^{(k)} = (d_1^k, d_2^k, \dots, d_n^k)$  表示之，稱之為第  $k$  階段的決定<sup>1)</sup>。矩陣  $P$  為  $\mathbf{d}^{(k)}$  的函數

$$P(\mathbf{d}^{(k)}) = \begin{pmatrix} p_{11}(d_1^{(k)}) & \cdots & p_{1n}(d_1^{(k)}) & p_{1,n+1}(d_1^{(k)}) \\ \cdots & \cdots & \cdots & \cdots \\ p_{n1}(d_n^{(k)}) & \cdots & p_{nn}(d_n^{(k)}) & p_{n,n+1}(d_n^{(k)}) \\ 0 & \cdots & 0 & 1 \end{pmatrix}$$

如果系統在時刻  $t_0$  之前所處的狀態給定，各時刻  $t_0, t_1, \dots$  前的決定  $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots$  都給定，則表征系統變化的隨機過程——馬爾可夫鏈也就確定。以  $\delta$  表示序列  $\{\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \dots\}$ ，並稱之為決策。於是，系統狀態變化過程的統計規律取決於所選取的決策。

以  $x_i$  表示由  $t_0$  前的初始狀態  $i (i = 1, 2, \dots, n)$  最終到達目標狀態的概率，則  $x_i$  為決策  $\delta$  的函數  $x_i(\delta)$ ；以  $\mathbf{x}(\delta)$  表示  $n$  維列向量  $\langle x_1(\delta), x_2(\delta), \dots, x_n(\delta) \rangle$ 。我們的問題是：

1. 是否存在這樣的  $\delta^*$ ，使得對一切其他的  $\delta$  都有： $\mathbf{x}(\delta^*) \geq \mathbf{x}(\delta)$ <sup>2)</sup>。 $\delta^*$  稱為最優決策。

2. 如果最優決策存在，怎樣求出  $\delta^*$  及相應的命中概率  $\mathbf{x}(\delta^*)$ 。

當各時刻  $t_0, t_1, \dots$  前所取的決定全都一樣時，決策  $\{\mathbf{d}, \mathbf{d}, \mathbf{d}, \dots\}$  稱為平穩決策。對於平穩決策，我們就以其決定元  $\mathbf{d}$  表示之。我們先證明，利用第三節的決策迭代法，從任意選取的平穩決策出發，經過有限次迭代，可以得到這樣的平穩決策  $\mathbf{d}^*$ ，使得對一切平穩決策  $\mathbf{d}$  都有  $\mathbf{x}(\mathbf{d}^*) \geq \mathbf{x}(\mathbf{d})$ 。然後證明，對於這一  $\mathbf{d}^*$  以及任何不平穩決策  $\delta$  都有  $\mathbf{x}(\mathbf{d}^*) \geq \mathbf{x}(\delta)$ 。於是，最優決策存在，可以在平穩決策中得到，而且可以用決策迭代法求出。

### 二、命中概率的性質

現在研究平穩決策下命中概率的性質。首先，在平穩決策下，所研究的过程形成一均勻馬爾可夫鏈<sup>3)</sup>，這鏈的統計規律由轉移概率矩陣  $P(\mathbf{d})$  完全確定。下面以矩陣  $P(\mathbf{d})$  表示這一馬爾可夫鏈，並稱之為鏈  $P(\mathbf{d})$ 。

根據馬爾可夫鏈的理論<sup>[10,11]</sup>，給定決策  $\mathbf{d}$  後，鏈  $P(\mathbf{d})$  的全部狀態可分成兩類：恆返狀態和非恆返狀態，而恆返狀態則組成若干恆返子鏈。顯然，目標狀態  $n + 1$  是恆返狀態，而且自身形成一恆返子鏈。如果鏈  $P(\mathbf{d})$  除了狀態  $n + 1$  所形成的恆返子鏈外，沒有其他的恆返子鏈，我們就稱決策  $\mathbf{d}$  為一致決策。

由命中概率的定義， $\mathbf{x}(\mathbf{d})$  應滿足等式

$$\mathbf{x}(\mathbf{d}) = Q(\mathbf{d})\mathbf{x}(\mathbf{d}) + \mathbf{q}(\mathbf{d}), \tag{1}$$

---

1)  $d_i^k$  表示控制的方式，不是數。但如果以號碼  $1, 2, \dots, l_i$  去標志各控制方式，則可以把  $d_i^k$  看作是這  $l_i$  個整數上取值的變數。本文將如此理解。  
 2) 本文以黑體拉丁字母  $\mathbf{a}$  表示分量為  $a_1, a_2, \dots, a_n$  的  $n$  維列向量，記作  $\mathbf{a} = \langle a_1, a_2, \dots, a_n \rangle$ 。記號  $\mathbf{a} \geq \mathbf{b}$  表示對所有  $i = 1, 2, \dots, n, a_i \geq b_i$ 。記號  $\mathbf{a} > \mathbf{b}$  表示對所有  $i = 1, 2, \dots, n, a_i \geq b_i$ ，同時至少有一  $i$  使得  $a_i > b_i$ 。  
 3) 即轉移概率不隨時間變化而保持固定的馬爾可夫鏈。

其中  $Q(\mathbf{d})$  是  $n \times n$  矩阵

$$Q(\mathbf{d}) = \begin{pmatrix} p_{11}(\mathbf{d}) & p_{12}(\mathbf{d}) & \cdots & p_{1n}(\mathbf{d}) \\ \cdots & \cdots & \cdots & \cdots \\ p_{n1}(\mathbf{d}) & p_{n2}(\mathbf{d}) & \cdots & p_{nn}(\mathbf{d}) \end{pmatrix}, \quad (2)$$

$\mathbf{q}(\mathbf{d})$  是  $n$  维列向量

$$\mathbf{q}(\mathbf{d}) = \langle p_{1,n+1}(\mathbf{d}), p_{2,n+1}(\mathbf{d}), \cdots, p_{n,n+1}(\mathbf{d}) \rangle. \quad (3)$$

考虑未知量  $\mathbf{x}$  的线性方程组

$$\mathbf{x} = Q(\mathbf{d})\mathbf{x} + \mathbf{q}(\mathbf{d}), \quad (4)$$

不难证明, 这方程组有唯一解的充分必要条件是决策  $\mathbf{d}$  为一致决策.

实际上, 如果  $\mathbf{d}$  不是一致决策, 则在状态  $1, 2, \cdots, n$  中存在有恒返子链. 这时, 矩阵  $I - Q(\mathbf{d})$  是降秩的, 因而(4)没有唯一解.

反之, 设  $\mathbf{d}$  是一致决策, 则状态  $1, 2, \cdots, n$  全是非恒返状态, 因之

对一切  $1 \leq i \leq n$ ,

$$\lim_{N \rightarrow \infty} p_{i,n+1}^{(N)}(\mathbf{d}) = 1^{(1)}.$$

故存在这样的  $N$ , 使得

对一切  $1 \leq i \leq n$ ,

$$\sum_{j=1}^n p_{ij}^{(N)}(\mathbf{d}) < 1,$$

因而矩阵  $Q(\mathbf{d})$  的最大特征数小于 1. 由此可知, 从任一初始向量  $\mathbf{x}^{(0)}$  开始, 利用递推公式

$$\mathbf{x}^{(k+1)} = Q(\mathbf{d})\mathbf{x}^{(k)} + \mathbf{q}(\mathbf{d}) \quad (5)$$

进行迭代, 序列  $\{\mathbf{x}^{(k)}\}$  必收敛于方程(4)的唯一解  $\mathbf{x} = \mathbf{1}^{(2)}$ .

当  $\mathbf{d}$  不是一致决策时, 状态  $1, 2, \cdots, n$  中含有恒返子链, 这些子链中任何状态的命中概率均为零, 即向量  $\mathbf{x}(\mathbf{d})$  中必含有零分量. 但这时  $\mathbf{x} = \mathbf{1}$  仍为方程组(4)的解, 可见它有无穷多个解. 在这种情况下, 从任意初始向量  $\mathbf{x}^{(0)}$  出发, 利用递推公式(5)进行迭代, 序列  $\{\mathbf{x}^{(k)}\}$  不一定收敛于  $\mathbf{x}(\mathbf{d})$ . 但是如果选择适当的初始向量  $\mathbf{x}^{(0)}$ , 则我们仍可得到这种收敛. 为此, 引入下述概念. 把除了目标状态而外的全部状态分成两类: 凡其命中概率大于零的统称优状态, 凡其命中概率等于零的统称劣状态. 当然, 谈到状态的优劣, 总是对给定的决策  $\mathbf{d}$  而言的.

如果我们这样选择初始向量  $\mathbf{x}^{(0)}$ , 使得对所有在决策  $\mathbf{d}$  下属于劣状态的状态  $i$ , 分量  $x_i^{(0)} = 0$ , 对属于优状态的状态  $i$ , 分量  $x_i^{(0)}$  可以任意选取 [下称此为条件(6)], 则由公式(5)所确定的序列  $\{\mathbf{x}^{(k)}\}$  必收敛于  $\mathbf{x}(\mathbf{d})$ :

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}(\mathbf{d}). \quad (6)$$

这一论断的证明见附录引理一.

1)  $p_{i,n+1}^{(N)}$  表示矩阵  $P$  的  $N$  次乘幂  $P^N$  的第  $i$  行第  $n+1$  列元素.

2)  $\mathbf{1}$  表示分量全是 1 的  $n$  维列向量.

### 三、決策迭代法

以  $\mathbf{d}^*$  表示平穩決策集合中的最優決策, 應用動態規劃的決策迭代法求  $\mathbf{d}^*$ , 這方法由下列兩個基本步驟組成。

**步驟 I:** 對前一次迭代所得的  $\mathbf{d}^{(r)}$  (在第一步迭代時, 則對任意選定的初始決策  $\mathbf{d}^{(0)}$ ), 計算相應的命中概率  $\mathbf{x}(\mathbf{d}^{(r)})$ 。為此, 只須列出矩陣  $P(\mathbf{d}^{(r)})$ , 令所有恆返子鏈中的狀態所對應的分量  $x_i(\mathbf{d}^{(r)}) = 0$ , 其餘非零分量則由方程組(4)中相應的等式聯立解出。

**步驟 II:** 根據前一步驟算出的  $\mathbf{x}(\mathbf{d}^{(r)})$ , 由下列公式

$$z_i^{(r+1)} = \max_{d_i} \left[ \sum_{j=1}^n p_{ij}(d_i) x_j(\mathbf{d}^{(r)}) + p_{i,n+1}(d_i) \right], \quad 1 \leq i \leq n \tag{7}$$

計算  $z_i^{(r+1)}$ 。當  $z_i^{(r+1)} = x_i(\mathbf{d}^{(r)})$  時, 保留  $d_i^{(r)}$  不變, 否則就用使上式方括弧取最大值的  $d_i$  替代  $d_i^{(r)}$ 。從  $i = 1$  到  $i = n$  進行這樣的計算和替代, 我們或者得到一個取代舊決策  $\mathbf{d}^{(r)}$  的新決策  $\mathbf{d}^{(r+1)}$ , 或者使決策  $\mathbf{d}^{(r)}$  不換。

如果將  $\mathbf{d}^{(r)}$  換成  $\mathbf{d}^{(r+1)}$ , 我們就繼續對  $\mathbf{d}^{(r+1)}$  施行步驟 I、II, 如此循環迭代, 直到決策不換為止。現在證明, 對任意選定的初始決策  $\mathbf{d}^{(0)}$ , 這種迭代程序能在有限步內收斂於最優平穩決策  $\mathbf{d}^*$ 。

首先注意, 不可能存在這樣的狀態: 在迭代前決策  $\mathbf{d}^{(r)}$  下是劣狀態, 而在迭代後決策  $\mathbf{d}^{(r+1)}$  下卻成為劣狀態。換句話說, 如果狀態  $i$  在決策  $\mathbf{d}^{(r+1)}$  下是劣狀態, 則有  $x_i(\mathbf{d}^{(r)}) = 0$ 。這事實從決策迭代法看來是直觀的, 但嚴格證明卻並非簡易的, 證明見附錄引理二。

考慮決策  $\mathbf{d}^{(r+1)}$  下的遞推公式(5), 並取  $\mathbf{x}(\mathbf{d}^{(r)})$  作為初始向量  $\mathbf{x}^{(0)}$ 。由上述事實, 對於決策  $\mathbf{d}^{(r+1)}$ ,  $\mathbf{x}^{(0)}$  滿足條件(6), 故由此確定的序列  $\{\mathbf{x}^{(k)}\}$  收斂:

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}(\mathbf{d}^{(r+1)}). \tag{8}$$

另一方面, 因

$$\begin{aligned} \mathbf{x}^{(1)} &= Q(\mathbf{d}^{(r+1)})\mathbf{x}^{(0)} + \mathbf{q}(\mathbf{d}^{(r+1)}) = Q(\mathbf{d}^{(r+1)})\mathbf{x}(\mathbf{d}^{(r)}) + \\ &\quad + \mathbf{q}(\mathbf{d}^{(r+1)}) = \mathbf{z}^{(r+1)} \geq \mathbf{x}(\mathbf{d}^{(r)}), \\ \mathbf{x}^{(2)} &= Q(\mathbf{d}^{(r+1)})\mathbf{x}^{(1)} + \mathbf{q}(\mathbf{d}^{(r+1)}) \geq Q(\mathbf{d}^{(r+1)})\mathbf{x}(\mathbf{d}^{(r)}) + \\ &\quad + \mathbf{q}(\mathbf{d}^{(r+1)}) = \mathbf{z}^{(r+1)} \geq \mathbf{x}(\mathbf{d}^{(r)}), \\ &\dots\dots\dots \end{aligned}$$

[依此類推, 對一切  $k = 1, 2, \dots$ , 都有  $\mathbf{x}^{(k)} \geq \mathbf{x}(\mathbf{d}^{(r)})$ ].

由(8)可得  $\mathbf{x}(\mathbf{d}^{(r+1)}) \geq \mathbf{x}(\mathbf{d}^{(r)})$ 。

不僅如此, 如果  $\mathbf{d}^{(r+1)} \neq \mathbf{d}^{(r)}$ , 則

$$\mathbf{x}(\mathbf{d}^{(r+1)}) > \mathbf{x}(\mathbf{d}^{(r)}).$$

事實上, 如果  $\mathbf{x}(\mathbf{d}^{(r+1)}) = \mathbf{x}(\mathbf{d}^{(r)})$ , 則這時

$$\begin{aligned} z_i^{(r+1)} &= \sum_{j=1}^n p_{ij}(d_i^{(r+1)}) x_j(\mathbf{d}^{(r)}) + p_{i,n+1}(d_i^{(r+1)}) = \\ &= \sum_{j=1}^n p_{ij}(d_i^{(r+1)}) x_j(\mathbf{d}^{(r+1)}) + p_{i,n+1}(d_i^{(r+1)}) = x_i(\mathbf{d}^{(r+1)}) = x_i(\mathbf{d}^{(r)}) \end{aligned}$$

1) 這裡  $\mathbf{z}^{(r+1)}$  為  $n$  維列向量, 其分量  $z_i^{(r+1)}$ , ( $1 \leq i \leq n$ ) 由(7)確定。

对  $i = 1, 2, \dots, n$  都成立, 因而  $\mathbf{d}^{(r+1)} = \mathbf{d}^{(r)}$  与假设矛盾.

由此可见, 每换一次决策, 指标必严格上升 (即至少有一分量加大), 所以不可能有同一决策循环出现的情况. 由于决策总数是有限的, 故在有限步内, 必出现决策不再替换的情况. 假设这时决策是  $\mathbf{d}^{(r)}$ , 则它必为最优平稳决策.

事实上, 这时对任一决策  $\mathbf{d}$  必有

$$\sum_{j=1}^n p_{ij}(\mathbf{d})x_j(\mathbf{d}^{(r)}) + p_{i,n+1}(\mathbf{d}) \leq x_i(\mathbf{d}^{(r)}), \quad 1 \leq i \leq n$$

即

$$Q(\mathbf{d})\mathbf{x}(\mathbf{d}^{(r)}) + \mathbf{q}(\mathbf{d}) \leq \mathbf{x}(\mathbf{d}^{(r)}).$$

考虑决策  $\mathbf{d}$  下的递推公式(5), 并取  $\mathbf{x}^{(0)} = \mathbf{0}^{(1)}$ , 则有

$$\begin{aligned} \mathbf{x}^{(1)} &= \mathbf{q}(\mathbf{d}) \leq Q(\mathbf{d})\mathbf{x}(\mathbf{d}^{(r)}) + \mathbf{q}(\mathbf{d}) \leq \mathbf{x}(\mathbf{d}^{(r)}), \\ \mathbf{x}^{(2)} &= Q(\mathbf{d})\mathbf{x}^{(1)} + \mathbf{q}(\mathbf{d}) \leq Q(\mathbf{d})\mathbf{x}(\mathbf{d}^{(r)}) + \mathbf{q}(\mathbf{d}) \leq \mathbf{x}(\mathbf{d}^{(r)}), \\ &\dots\dots\dots \end{aligned}$$

依此类推, 对一切  $k = 1, 2, \dots$ , 都有

$$\mathbf{x}^{(k)} \leq \mathbf{x}(\mathbf{d}^{(r)}).$$

但

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}(\mathbf{d}),$$

故

$$\mathbf{x}(\mathbf{d}) \leq \mathbf{x}(\mathbf{d}^{(r)}). \quad (\text{对一切 } \mathbf{d})$$

因此  $\mathbf{d}^{(r)}$  是最优平稳决策  $\mathbf{d}^*$ . 与此同时, 我们就证明了在平稳决策集中最优决策是存在的<sup>2)</sup>.

### 四、考虑不平稳决策的情况

现在要证明, 上面所求得的最优平稳决策  $\mathbf{d}^*$  在包含不平稳决策的决策集中也是最优的. 为此, 先推导最优指标函数  $\mathbf{x}(\mathbf{d}^*)$  的一个重要性质.

以  $\mathbf{x}^*$  记  $\mathbf{x}(\mathbf{d}^*)$ , 根据动态规划的最优化原理<sup>3)</sup>,  $\mathbf{x}^*$  应满足下列方程

$$\mathbf{x}^* = \max_{\mathbf{d}} [Q(\mathbf{d})\mathbf{x}^* + \mathbf{q}(\mathbf{d})]. \tag{9}$$

下面证明, 这方程的解  $\mathbf{x}^*$  可以由动态规划的函数迭代法求出. 就是说, 如果取  $\mathbf{x}^{(0)} = \mathbf{0}$ , 并依下列递推公式定义序列  $\{\mathbf{x}^{(k)}\}$ :

$$\mathbf{x}^{(k+1)} = \max_{\mathbf{d}} [Q(\mathbf{d})\mathbf{x}^{(k)} + \mathbf{q}(\mathbf{d})], \quad k = 0, 1, 2, \dots \tag{10}$$

则有

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*. \tag{11}$$

为了简化起见, 我们把式(10)写成

$$\mathbf{x}^{(k+1)} = T\mathbf{x}^{(k)},$$

考虑决策  $\mathbf{d}^*$  下的递推公式(5)

$$\mathbf{x}^{(k+1)} = Q(\mathbf{d}^*)\mathbf{x}^{(k)} + \mathbf{q}(\mathbf{d}^*),$$

1)  $\mathbf{0}$  表示分量全是零的  $n$  维列向量.  
2) 对于最优平稳决策的存在性, 可以给出另一完全独立的证明.  
3) 也可以由第三节  $\mathbf{d}^*$  的求法得出.

或者写成

$$\mathbf{x}^{(k+1)} = T^* \mathbf{x}^{(k)}.$$

因为对任意的  $\mathbf{x} \geq \mathbf{0}$ ,

$$\max_{\mathbf{d}} [Q(\mathbf{d})\mathbf{x} + \mathbf{q}(\mathbf{d})] \geq Q(\mathbf{d}^*)\mathbf{x} + \mathbf{q}(\mathbf{d}^*),$$

即

$$T\mathbf{x} \geq T^*\mathbf{x};$$

当  $\mathbf{x} = \mathbf{x}^{(0)} = \mathbf{0}$  时,得

$$T\mathbf{x}^{(0)} \geq T^*\mathbf{x}^{(0)}.$$

由这两个不等式推得

$$T^{(2)}\mathbf{x}^{(0)} = T(T\mathbf{x}^{(0)}) \geq T^*(T^*\mathbf{x}^{(0)}) = T^{*(2)}\mathbf{x}^{(0)}.$$

依此,对一切  $k = 1, 2, \dots$ , 可推得

$$T^{(k)}\mathbf{x}^{(0)} \geq T^{*(k)}\mathbf{x}^{(0)}. \tag{12}$$

因  $\mathbf{x}^{(0)} = \mathbf{0}$ , 不論  $\mathbf{d}^*$  如何,条件(6)恆满足,故

$$\lim_{k \rightarrow \infty} T^{*(k)}\mathbf{x}^{(0)} = \mathbf{x}^*.$$

于是,由式(12)得

$$\liminf_{k \rightarrow \infty} T^{(k)}\mathbf{x}^{(0)} \geq \mathbf{x}^*. \tag{13}$$

另一方面,注意到  $\mathbf{x}^* \geq \mathbf{0}$ , 而  $\mathbf{x}^*$  满足式(9),故对任意平穩决策  $\mathbf{d}$ , 我們有

$$T\mathbf{x}^{(0)} = \max_{\mathbf{d}} [Q(\mathbf{d})] \leq \max_{\mathbf{d}} [Q(\mathbf{d})\mathbf{x}^* + \mathbf{q}(\mathbf{d})] = \mathbf{x}^*.$$

依此,对一切  $k = 1, 2, \dots$  可推得

$$T^{(k)}\mathbf{x}^{(0)} \leq \mathbf{x}^*.$$

于是,

$$\limsup_{k \rightarrow \infty} T^{(k)}\mathbf{x}^{(0)} \leq \mathbf{x}^*. \tag{14}$$

合并式(13)和(14),就得到式(11).

現在考虑任一不平穩决策  $\delta$ , 而  $\delta = \{\mathbf{d}^{(0)}, \mathbf{d}^{(1)}, \mathbf{d}^{(2)}, \dots\}$ . 令

$$\begin{aligned} \mathbf{y}^{(1)} &= \mathbf{q}(\mathbf{d}^{(0)}), \\ \mathbf{y}^{(2)} &= \mathbf{y}^{(1)} + Q(\mathbf{d}^{(0)})\mathbf{q}(\mathbf{d}^{(1)}), \\ \mathbf{y}^{(3)} &= \mathbf{y}^{(2)} + Q(\mathbf{d}^{(0)})Q(\mathbf{d}^{(1)})\mathbf{q}(\mathbf{d}^{(2)}), \\ &\dots\dots\dots \\ \mathbf{y}^{(k)} &= \mathbf{y}^{(k-1)} + Q(\mathbf{d}^{(0)})Q(\mathbf{d}^{(1)})\dots Q(\mathbf{d}^{(k-2)})\mathbf{q}(\mathbf{d}^{(k-1)}), \\ &\dots\dots\dots \end{aligned}$$

由命中概率的定义可得,在决策  $\delta$  下的命中概率  $\mathbf{x}(\delta)$  为

$$\mathbf{x}(\delta) = \lim_{k \rightarrow \infty} \mathbf{y}^{(k)}. \tag{15}$$

把序列  $\{\mathbf{y}^{(k)}\}$  与由式(10)所确定的且初始向量  $\mathbf{x}^{(0)} = \mathbf{0}$  的序列  $\{\mathbf{x}^{(k)}\}$  进行比较:

$$\begin{aligned} \mathbf{x}^{(1)} &= \max_{\mathbf{d}} [Q(\mathbf{d})] \geq \mathbf{q}(\mathbf{d}^{(0)}) = \mathbf{y}^{(1)}, \\ \mathbf{x}^{(2)} &= \max_{\mathbf{d}} [Q(\mathbf{d}) + Q(\mathbf{d})\mathbf{x}^{(1)}] \geq \mathbf{q}(\mathbf{d}^{(0)}) + Q(\mathbf{d}^{(0)})\mathbf{x}^{(1)} \geq \\ &\geq \mathbf{y}^{(1)} + Q(\mathbf{d}^{(0)})\mathbf{q}(\mathbf{d}^{(1)}) = \mathbf{y}^{(2)}, \\ &\dots\dots\dots \end{aligned}$$



(A.3) 右边  $x_j$  的系数矩陣的最大特征数小于 1。由此可知, 从任何初始的  $x_{s+1}^{(0)}, x_{s+2}^{(0)}, \dots, x_n^{(0)}$  出发, 利用公式 (5) 进行迭代, 所得  $x_{s+1}^{(k)}, x_{s+2}^{(k)}, \dots, x_n^{(k)}$  必收敛于方程組 (A.3) 的唯一解, 但  $x_{s+1}(\mathbf{d}), x_{s+2}(\mathbf{d}), \dots, x_n(\mathbf{d})$  是它的一个解, 故得

$$\lim_{k \rightarrow \infty} x_i^{(k)} = x_i(\mathbf{d}), \quad s+1 \leq i \leq n \tag{A.4}$$

合并式(A.2)和式(A.4), 就得式(A.1)。

**引理二** 設  $\mathbf{d}^{(r)}$  是迭代前决策,  $\mathbf{d}^{(r+1)}$  是迭代后决策, 則  $G(\mathbf{d}^{(r)}) \cap B(\mathbf{d}^{(r+1)})$  是空集。

**証明** 为簡化起見, 引入下列記号:  $GG = G(\mathbf{d}^{(r)}) \cap G(\mathbf{d}^{(r+1)}), GB = G(\mathbf{d}^{(r)}) \cap B(\mathbf{d}^{(r+1)}), BG = B(\mathbf{d}^{(r)}) \cap G(\mathbf{d}^{(r+1)}), GG = G(\mathbf{d}^{(r)}) \cap G(\mathbf{d}^{(r+1)})$ 。

采用反証法。假定  $GB$  不是空集。因对  $i \in GB, j \in GG \cup BG, p_{ij}(d_i^{(r+1)}) = 0$ , 而对  $j \in BB, x_j(\mathbf{d}^{(r)}) = 0$ , 故  $x_i^{(r+1)} = \sum_{j \in GB} p_{ij}(d_i^{(r+1)}) x_j(\mathbf{d}^{(r)}) \geq x_i(\mathbf{d}^{(r)})$ 。显然, 这时至少存在  $i \in GB$ , 使得

$$x_i^{(r+1)} = \sum_{j \in GB} p_{ij}(d_i^{(r+1)}) x_j(\mathbf{d}^{(r)}) = x_i(\mathbf{d}^{(r)})$$

以  $\overline{GB}$  記  $GB$  中使上述等式成立之子集,  $\widehat{GB}$  記  $\overline{GB}$  对  $GB$  之余集, 則  $\overline{GB}$  不是空集。这时可能出现两种情况。其一,  $\widehat{GB}$  是空集, 其二,  $\widehat{GB}$  不空。

先考虑第一种情况。这时,  $GB = \overline{GB}$ , 故由决策迭代規則, 对所有  $i \in GB, d_i^{(r+1)} = d_i^{(r)}$ , 于是  $p_{ij}(d_i^{(r+1)}) = p_{ij}(d_i^{(r)})$ 。因而对所有  $i \in GB, j \in GG \cup BG$  及  $j = n+1, p_{ij}(d_i^{(r)}) = p_{ij}(d_i^{(r+1)}) = 0$ 。考虑决策  $\mathbf{d}^{(r)}$  下的递推公式(5), 当  $i \in GB$  时,

$$x_i^{(k+1)} = \sum_{j \in GB \cup BB} p_{ij}(d_i^{(r)}) x_j^{(k)} \tag{A.5}$$

因为对  $i \in BB \cup BG, j = GB \cup GG$  及  $j = n+1, p_{ij}(d_i^{(r)}) = 0$ , 故对  $i \in BB \cup BG,$

$$x_i^{(k+1)} = \sum_{j \in BB \cup BG} p_{ij}(d_i^{(r)}) x_j^{(k)} \tag{A.6}$$

因此, 如果取  $\mathbf{x}^{(0)}$  使得对所有  $i \in GB \cup BB \cup BG, x_i^{(0)} = 0$ , 則由式(A.5), (A.6) 可見, 对这些  $i,$

$$x_i^{(k)} = 0, \quad k = 1, 2, \dots$$

另一方面, 由引理一,  $\lim_{k \rightarrow \infty} x_i^{(k)} = x_i(\mathbf{d}^{(r)})$ 。于是对  $i \in GB$  得  $x_i(\mathbf{d}^{(r)}) = 0$ , 这与  $GB$  的定义矛盾。

其次考虑第二种情况。如果  $\max_{i \in \overline{GB}} x_j(\mathbf{d}^{(r)}) > \max_{j \in \widehat{GB}} x_j(\mathbf{d}^{(r)})$ , 則因对  $i \in \overline{GB}, d_i^{(r+1)} = d_i^{(r)}$ , 故

$$x_i(\mathbf{d}^{(r)}) = \sum_{j \in \overline{GB}} p_{ij}(d_i^{(r)}) x_j(\mathbf{d}^{(r)}) + \sum_{j \in \widehat{GB}} p_{ij}(d_i^{(r)}) x_j(\mathbf{d}^{(r)})$$

可見在决策  $\mathbf{d}^{(r)}$  下,  $\overline{GB}$  中使  $x_j(\mathbf{d}^{(r)})$  达到  $\max_{j \in \overline{GB}} x_j(\mathbf{d}^{(r)})$  的那些状态組成一恒返子鏈, 因而对这些状态,

$x_j(\mathbf{d}^{(r)}) = 0$ , 这和  $\overline{GB}$  的定义相矛盾。反之, 如果  $\max_{j \in \overline{GB}} x_j(\mathbf{d}^{(r)}) \leq \max_{j \in \widehat{GB}} x_j(\mathbf{d}^{(r)})$ , 則对  $i \in \widehat{GB},$

$$x_i^{(r+1)} = \sum_{j \in \overline{GB}} p_{ij}(d_i^{(r+1)}) x_j(\mathbf{d}^{(r)}) + \sum_{j \in \widehat{GB}} p_{ij}(d_i^{(r+1)}) x_j(\mathbf{d}^{(r)})$$

$$x_j(\mathbf{d}^{(r)}) \leq \max_{j \in \widehat{GB}} x_j(\mathbf{d}^{(r)}) \sum_{i \in \overline{GB}} p_{ij}(d_i^{(r+1)}) \leq \max_{j \in \widehat{GB}} x_j(\mathbf{d}^{(r)})$$

这又与  $\widehat{GB}$  的定义相矛盾。因此,  $GB$  不是空集是不可能的。

### 参 考 文 献

- [1] Bellman, R., Dynamic Programming, Princeton University Press, Princeton, New Jersey, 1957.
- [2] Bellman, R., A Markovian Decision Process, Journal of Mathematics and Mechanics, 6 (1957), No. 5, 679—684.
- [3] Howard, R. A., Dynamic Programming and Markov Processes, John Wiley & Sons Inc., New York, 1960 (动态规划与馬爾柯夫过程, 李为政等譯, 上海科学技术出版社, 1963)。

- [4] Eaton, J. H. & Zadeh, L. A., Optimal Pursuit Strategies in Discrete-state Probabilistic Systems, *Trans. of ASME, Journal of Basic Engineering*, **84** (1962), No. 1, 23—29.
- [5] Blackwell, D. V., On the Functional Equation of Dynamic Programming, *Journal of Mathematical Analysis and Applications*, **2** (1961), No. 2, 273—276.
- [6] Blackwell, D. V., Discrete dynamic Programming, *The Annals of Mathematical Statistics*, **33** (1962), No. 2, 719—726.
- [7] Derman, C., On Sequential Decisions and Markov Chains, *Management Science*, **9** (1962), No. 1, 16—24.
- [8] Derman, C., Stable Sequential Control Rules and Markov Chains, *Journal of Mathematical Analysis and Applications*, **6** (1963), No. 2, 257—265.
- [9] Понтрягин, Л. С., Болтянский, В. Г., Гамкрелидзе, Р. В., Мищенко, Е. Ф., Математическая теория оптимальных процессов, ГИФМЛ, Москва, 1961.
- [10] Feller, W., An Introduction to Probability Theory and its Applications, John Wiley & Sons Inc., 1957.
- [11] Романовский, В. И., Дискретные цепи Маркова, ГИТТЛ, Москва, 1949 (疏散的馬爾可夫鏈, 梁友琪譯, 科学出版社, 北京, 1958).

## AN OPTIMAL POLICY FOR CONTROLLING THE CONTROLLABLE MARKOV CHAINS

WU CHANG-PU

This paper is concerned with one type of the optimal Markov controlled systems. The controlled system is described by a Markov chain whose statistical property depends on the sequence of decisions that we call a policy. There exists an objective state with the property that once the system reaches this state, it remains unchanged forever. Our purpose is to choose a policy which maximizes all the probabilities that the system ever reaches this objective state from every initial state. First we give a policy-iteration method for obtaining an optimal policy over the set of stable policies. We then prove such a policy is also optimal over the set containing both stable and unstable policies.