

Управление марковским процессом в задаче с ограничениями

В настоящей работе рассматривается условная задача динамического программирования с затухающим действием в случае любого конечного числа ограничений и конечных множеств состояний и управлений. Безусловная задача хорошо известна.

Условную задачу с одним ограничением рассмотрел Фрид [1]. К сожалению, подход к задаче, использованный в его работе, остался незамеченным из-за большого числа чисто технических приемов. В этой работе, являющейся продолжением [1], ставится цель не только обобщить задачу на случай произвольного конечного числа ограничений, но и выделить основную идею, не затемняя ее второстепенными выкладками. Поэтому рассматриваются лишь конечные множества состояний и управлений, хотя данный подход позволяет получать аналогичные результаты и в более общих случаях (см., например, [1]).

Подход к решению этой задачи, как и в [1], основан на принципе Лагранжа, что позволяет получить необходимые и достаточные условия оптимальности стратегии, а также выделить определенный класс стратегий, в котором существует хотя бы одна оптимальная. Основным результатом этой работы заключается в следующем: существует оптимальная стратегия, являющаяся взвесью не более чем $m + 1$ -й стационарной стратегии, где m — число ограничений.

Пусть S и A — непустые конечные множества. Пусть f_0, f_1, \dots, f_m — ограниченные функции на $S \times A$, а q ставит в соответствие каждой паре $(s, a) \in S \times A$ вероятностную меру $q(\cdot | s, a)$, определенную на множестве состояний S , в соответствии с которой происходит переход из состояния s_i в s_{i+1} . Положим $\mathcal{H}_n = S_1 \times A_1 \times \dots \times S_n$, где $S_i \equiv S$, $A_i \equiv A$ для любого $i = 1, \dots, n$. Стратегией π будем называть последовательность функций $\{\pi^n(\cdot | h_n)\}$, $n = 1, 2, \dots$, $h_n \in \mathcal{H}_n$, каждая из которых является переходной функцией из \mathcal{H}_n в A . Класс всех возможных стратегий обозначим через Π . Введем понятие взвешенной стратегии.

Определение. α^n -взвесью стратегий π_1, \dots, π_n с весом $\alpha^n = (\alpha_1, \dots, \alpha_n)$, где $\sum_{i=1}^n \alpha_i = 1$, $\alpha_i \geq 0$, назовем стратегию $\pi(\alpha^n, \pi_1, \dots, \pi_n)$, которую определим как последовательность переходных функций $\{\pi^l(\alpha^n, \pi_1, \dots, \pi_n)\}$, $l = 1, 2, \dots$, каждая из которых $\pi^l(\alpha^n, \pi_1, \dots, \pi_n) = \sum_{i=1}^n \alpha_i \pi_i^l(\cdot | h_i)$.

Таким образом, Π — выпуклое множество.

Для любого $n \geq 2$ и $s_1 \in S$ распределение q и стратегия π определяют меру на \mathcal{H}_n . В силу теоремы Ионеску — Тулча [2] при фиксированном s_1 распределение q и стратегия π определяют меру $P_{s_1}^\pi$ на пространстве траекторий $S_1 \times A_1 \times S_2 \times A_2 \times \dots$, которое обозначим Ω . Математическое ожидание по этой мере будем обозначать $M_{s_1}^\pi$. В дальнейшем, полагая s_1 всегда фиксированным, будем опускать этот индекс.

Замечание 1. Пусть F — ограниченная функция на пространстве траекторий Ω и n стратегий π_1, \dots, π_n , тогда $M^\lambda F = \sum_{i=1}^n \alpha_i M^{\pi_i} F$, где $\lambda = \pi(\alpha^n, \pi_1, \dots, \pi_n)$.

Рандомизированной марковской стратегией будем называть такую стратегию, в которой π^n зависит только от состояния s_n , в котором система находится в текущий момент времени. Марковская стратегия — последователь-

ность $\{g_i\}$, $i = 1, 2, \dots$, где каждая g_i — измеримая функция, отображающая S в A , и $g_n(s)$ — управление, которое выбираем на n -м шагу, если n -е состояние — s . Стационарная стратегия — марковская стратегия, в которой $g_n = g$ для всех n , где g — некоторая измеримая функция, отображающая S в A . Понятно, что класс всех стационарных стратегий в нашем случае (S и A конечны) конечен.

Рассмотрим $\beta \in [0, 1)$. Определим функции F_i^β , $i = 0, \dots, m$, на пространстве траекторий Ω : $F_i^\beta = \sum_{j=1}^{\infty} \beta^{j-1} f_i(s_j, a_j)$. Далее будем опускать индекс β , полагая его везде фиксированным. Заметим, что для любого фиксированного $\beta \in [0, 1)$ функции F_i ограничены.

Зафиксируем $K = (K_1, \dots, K_m)$, где $K_i < \infty$ для любого $i = 1, \dots, m$.

Будем говорить, что стратегия π принадлежит классу \mathcal{D} , если $M^\pi F_i \leq K_i$ для любого $i = 1, \dots, m$. Предположим, что класс \mathcal{D} непуст, и назовем ценой игры в точке s_1 число $v = \sup_{\pi \in \mathcal{D}} M^\pi F_0$. Стратегию $\pi_0 \in \mathcal{D}$, на которой достигается v (если она существует), назовем оптимальной.

Поставим перед собой задачу исследовать существование оптимальной стратегии, а также выяснить, к какому классу стратегий могла бы принадлежать π_0 . Существенно, что в отличие от безусловной задачи, где доказывалось существование оптимальной стратегии сразу во всех точках, в условной задаче может не существовать стратегии, оптимальной даже в двух точках (см. пример 2 из [1]).

Если зафиксировать функцию F на Ω , то функция $M^\pi F$ как функция от стратегии π является линейной и ограниченной, если ограничена функция F . Как уже было сказано, множество всех стратегий Π — выпуклое множество. Поэтому поставленная задача, которую формально можно представить в виде

$$M^\pi F_0 \rightarrow \sup, \quad (1)$$

$$M^\pi F_i - K_i \leq 0, \quad i = 1, \dots, m, \quad (2)$$

$$\pi \in \Pi, \quad (3)$$

является выпуклой экстремальной задачей оптимального управления.

Введем функцию Лагранжа $\mathcal{L}(\pi, \lambda, \lambda_0) = \lambda_0 M^\pi F_0 - \sum_{i=1}^m \lambda_i (M^\pi F_i - K_i)$, где $\lambda = (\lambda_1, \dots, \lambda_m)$. Согласно теореме Куна — Таккера [3] справедливо:

1) если π_0 — оптимальная стратегия, то существует вектор $(\hat{\lambda}_0, \hat{\lambda}_1, \dots, \hat{\lambda}_m) \neq 0$ такой, что выполняются условия

$$\hat{\lambda}_i \geq 0, \quad i = 0, \dots, m, \quad (4)$$

$$\hat{\lambda}_i (M^{\pi_0} F_i - K_i) = 0, \quad i = 1, \dots, m, \quad (5)$$

$$\mathcal{L}(\pi_0, \hat{\lambda}, \hat{\lambda}_0) = \sup_{\pi \in \Pi} \mathcal{L}(\pi, \hat{\lambda}, \hat{\lambda}_0); \quad (6)$$

2) если $\hat{\lambda}_0 > 0$ и $\pi_0 \in D$, то условия (4) — (6) являются достаточными для того, чтобы стратегия π_0 была оптимальной;

3) для того чтобы $\hat{\lambda}_0 \neq 0$ достаточно, чтобы нашлась стратегия $\bar{\pi} \in \Pi$ такая, что выполняются условия: $M^{\bar{\pi}} F_i - K_i < 0$, $i = 1, \dots, m$.

Покажем, что найдется вектор $(1, \lambda) \in \mathbb{R}_+^{m+1}$ и стратегия $\pi_0 \in \mathcal{D}$ такие, что выполняются условия (4) — (6). Как обычно $\mathbb{R}_+^{m+1} = \{x \in \mathbb{R}^{m+1}; x_i \geq 0, i = 0, \dots, m\}$. Тем самым докажем существование оптимальной стратегии и укажем класс стратегий, в котором она всегда существует.

Введем функцию $\omega(\lambda) = \sup_{\pi \in \Pi} \mathcal{L}(\pi, \lambda, 1)$. Так как мы ищем вектор

$(1, \hat{\lambda}) \in \mathbb{R}_+^{m+1}$, то условие (6) принимает вид

$$\omega(\hat{\lambda}) = \sup_{\pi \in \Pi} \mathcal{L}(\pi, \hat{\lambda}, 1) = \mathcal{L}(\pi_0, \hat{\lambda}, 1). \quad (7)$$

Функция $\mathcal{L}(\pi, \lambda, 1)$ линейна по λ при любой фиксированной стратегии π . Поскольку справедлива цепочка равенств

$$\begin{aligned} F_0 - \sum_{i=1}^m \lambda_i (F_i - K_i) &= \sum_{j=1}^{\infty} \beta^{j-1} f_0(s_j, a_j) - \sum_{i=1}^m \lambda_i \left(\sum_{j=1}^{\infty} \beta^{j-1} f_i(s_j, a_j) - K_i \right) = \\ &= \sum_{j=1}^{\infty} \beta^{j-1} (f_0(s_j, a_j) - \sum_{i=1}^m \lambda_i f_i(s_j, a_j)) + \sum_{i=1}^m \lambda_i K_i, \end{aligned}$$

функция $f_0(s, a) - \sum_{i=1}^m \lambda_i f_i(s, a)$ ограничена, а сумма $\sum_{i=1}^m \lambda_i K_i = \text{const}$ при любом фиксированном λ , то согласно теореме 7 (b) из [4] точная верхняя грань функции $\mathcal{L}(\pi, \lambda, 1)$ по π достигается в любой фиксированной точке λ на стационарной стратегии.

Поэтому для исследования функции $\omega(\lambda)$ достаточно рассматривать функцию $\mathcal{L}(\pi, \lambda, 1)$ только на таких стратегиях. Как было замечено, в случае конечных множеств состояний и управлений стационарных стратегий конечно число, поэтому справедливо:

1) для любого фиксированного λ существует стационарная стратегия $\pi(\lambda)$ на которой достигается $\omega(\lambda)$. Заметим, что этой фразой стратегия $\pi(\lambda)$ определяется, вообще говоря, неоднозначно;

2) $\omega(\lambda)$ как функция от λ представляет собой выпуклую вниз кусочно-линейную функцию, причем число таких «кусков» не превышает числа стационарных стратегий, а значит конечно.

Рассмотрим в \mathbb{R}^{m+1} график функции $\omega(\lambda)$, который обозначим W . Графиком функции $\mathcal{L}(\pi, \lambda, 1)$ для каждой фиксированной стратегии π будет гиперплоскость, которую обозначим $T(\pi)$. Будем рассматривать также гиперплоскости $T(c) = \{(c, \lambda) : \lambda \in \mathbb{R}^m, c = \text{const}\}$.

Для того чтобы лучше понять смысл дальнейших рассуждений, предположим, что имеем оптимальную стратегию π_0 в задаче (1)–(3) и точку $\hat{\lambda} \in \mathbb{R}_+^m$ такую, что $\hat{\lambda}$ удовлетворяет (5) и $\hat{\lambda}_i \neq 0$ для всех $i = 1, \dots, m$. Тогда $M^{\pi_0} F_i - K_i = 0$, $i = 1, \dots, m$, или, другими словами, гиперплоскость $T(\pi_0)$ совпадает с $T(\hat{c})$, если $\hat{c} = M^{\pi_0} F_0$.

Пусть теперь часть координат точки $\hat{\lambda}$ равны нулю. Тогда для тех номеров i , при которых $\hat{\lambda}_i \neq 0$, из условия (5) следует строгое равенство $M^{\pi_0} F_i = K_i$, а для остальных i выполняются неравенства $M^{\pi_0} F_i - K_i \leq 0$, так как $\pi_0 \in \mathcal{D}$. Но тогда $\omega(\lambda) \geq \mathcal{L}(\pi_0, \lambda, 1) = M^{\pi_0} F_0 - \sum_{i=1}^m \lambda_i (M^{\pi_0} F_i - K_i) \geq M^{\pi_0} F_0$ для любого $\lambda \in \mathbb{R}_+^m$. Положим $\hat{c} = M^{\pi_0} F_0$, тогда $\omega(\hat{\lambda}) = \mathcal{L}(\pi_0, \hat{\lambda}, 1) = M^{\pi_0} F_0 = \hat{c}$, т. е. $\omega(\hat{\lambda}) = \hat{c}$ и $\omega(\lambda) \geq \hat{c}$ для любого $\lambda \in \mathbb{R}_+^m$. Другими словами, график W лежит не ниже гиперплоскости $T(\hat{c})$ при $\lambda \in \mathbb{R}_+^m$ и у них существует общая точка $(\hat{c}, \hat{\lambda})$.

Будем искать точку $\hat{\lambda}$, надвигая снизу гиперплоскость $T(c)$ на график W , рассматриваемый только в \mathbb{R}_+^m . Возникает вопрос, всегда ли можно это сделать? В лемме 1 покажем, что \hat{c} и точка «первого касания» $\lambda^* \in \mathbb{R}_+^m$ существуют всегда, если \mathcal{D} непусто, а в леммах 2 и 3 выясним, что точку λ^* можно взять за искомую $\hat{\lambda}$ и укажем оптимальную стратегию π_0 .

Лемма 1. Если множество \mathcal{D} непусто, то существует константа \hat{c} и точка $\lambda^* \in \mathbb{R}_+^m$ такие, что $\hat{c} = \omega(\lambda^*) = \inf_{\lambda \in \mathbb{R}_+^m} \omega(\lambda) > -\infty$.

Доказательство. Так как множество \mathcal{D} непусто, то существует стратегия π такая, что $M^{\pi}F_i - K_i \leq 0, i = 1, \dots, m$. Положим $c = M^{\pi}F_0$, тогда $c = M^{\pi}F_0 \leq M^{\pi}F_0 - \sum_{i=1}^m \lambda_i (M^{\pi}F_i - K_i) \leq \omega(\lambda)$ при любом $\lambda \in \mathbb{R}_+^m$, т. е. функция $\omega(\lambda)$ ограничена снизу на \mathbb{R}_+^m . Поскольку она кусочно линейная и таких «кусочков» конечное число, существует точка λ^* такая, что $\hat{c} = \omega(\lambda^*) = \inf_{\lambda \in \mathbb{R}_+^m} \omega(\lambda) > -\infty$. Лемма доказана.

Замечание 2. Так как $\omega(\lambda)$ — выпуклая функция, то существуют односторонние частные производные. Нетрудно заметить, что в точке λ^* , если $\lambda_i^* \neq 0$, должно выполняться неравенство $\omega'_{\lambda_i^-}(\lambda^*) \leq 0 \leq \omega'_{\lambda_i^+}(\lambda^*)$.

Рассмотрим все стратегии $\pi(\lambda^*)$ и пронумеруем их в произвольном порядке. Обозначим $\pi_i(\lambda^*)$ через $\pi^*(i)$. Их конечное число, которое обозначим через k . Заметим, что любой стратегии $\pi = \pi(\alpha^1, \pi_1, \dots, \pi_n)$ соответствует функция $\mathcal{L}(\pi, \lambda, 1) = \sum_{i=1}^n \alpha_i \mathcal{L}(\pi_i, \lambda, 1)$, т. е. являющаяся линейной комбинацией функций $\mathcal{L}(\pi_i, \lambda, 1)$ с теми же весами α_i .

Лемма 2. Если найдется вектор $\hat{\alpha} = (\hat{\alpha}_1, \dots, \hat{\alpha}_k)$ такой, что $\sum_{i=1}^k \hat{\alpha}_i = 1$ и $\hat{\alpha}_i \geq 0, i = 1, \dots, k$, и при этом будет выполнено неравенство $\sum_{i=1}^k \hat{\alpha}_i \mathcal{L}(\pi^*(i), \lambda, 1) \geq \hat{c}$ для любого $\lambda \in \mathbb{R}_+^m$, где $\hat{c} = \inf_{\lambda \in \mathbb{R}_+^m} \omega(\lambda)$, то стратегия $\pi_0 = \pi(\hat{\alpha}, \pi^*(1), \dots, \pi^*(k))$ — оптимальная.

Доказательство. Проверим достаточные условия оптимальности стратегии.

Условие (4) и условие $\lambda_0^* > 0$ выполнены, поскольку рассматриваем вектор $(1, \lambda^*), \lambda^* \in \mathbb{R}_+^m$.

Условие (6) эквивалентно условию (7), которое выполняется, так как $\mathcal{L}(\pi_0, \lambda^*, 1) = \sum_{i=1}^k \hat{\alpha}_i \mathcal{L}(\pi^*(i), \lambda^*, 1) = \sum_{i=1}^k \hat{\alpha}_i \omega(\lambda^*) = \omega(\lambda^*) = \sup_{\pi \in \Pi} \mathcal{L}(\pi, \lambda^*, 1)$.

Проверим условие (5). Пусть сначала найдется номер i такой, что λ_i^* строго больше нуля. Рассмотрим точки $\lambda_1 = (\lambda_1^*, \dots, \lambda_{i-1}^*, \lambda_i^* - \lambda_0, \lambda_{i+1}^*, \dots, \lambda_m^*), \lambda_2 = (\lambda_1^*, \dots, \lambda_{i-1}^*, \lambda_i^* + \lambda_0, \lambda_{i+1}^*, \dots, \lambda_m^*)$, где $\lambda_0 > 0, \lambda_i^* - \lambda_0 > 0$.

Тогда $\hat{c} \leq \sum_{i=1}^k \hat{\alpha}_i \mathcal{L}(\pi^*(i), \lambda_1, 1) = \mathcal{L}(\pi_0, \lambda_1, 1) = \mathcal{L}(\pi_0, \lambda^*, 1) + \lambda_0 (M^{\pi_0}F_i - K_i) = \hat{c} + \lambda_0 (M^{\pi_0}F_i - K_i)$, т. е.

$$0 \leq \lambda_0 (M^{\pi_0}F_i - K_i), \quad (8)$$

$$\begin{aligned} \hat{c} &\leq \sum_{i=1}^k \hat{\alpha}_i \mathcal{L}(\pi^*(i), \lambda_2, 1) = \mathcal{L}(\pi_0, \lambda_2, 1) = \mathcal{L}(\pi_0, \lambda^*, 1) - \lambda_0 (M^{\pi_0}F_i - K_i) = \\ &= \hat{c} - \lambda_0 (M^{\pi_0}F_i - K_i), \end{aligned}$$

т. е.

$$0 \geq \lambda_0 (M^{\pi_0}F_i - K_i). \quad (9)$$

Из (8) и (9) заключаем, что $M^{\pi_0}F_i - K_i = 0$ для любого номера i такого, что $\lambda_i^* > 0$. Если $\lambda_i^* = 0$, то условие (5) выполняется для этого номера i очевидным образом.

Теперь проверим, что $\pi_0 \in \mathcal{D}$.

Для этого осталось проверить, что $M^{\pi^*} F_i - K_i \leq 0$ для всех номеров i таких, что $\lambda_i^* = 0$. Но это также следует из неравенства (9). Лемма доказана.

Рассмотрим функцию $\mathcal{L}(\lambda) = \max_{1 \leq i \leq k} \mathcal{L}(\pi^*(i), \lambda, 1)$. Очевидно, что $\mathcal{L}(\lambda^*) = \hat{c}$ и $\mathcal{L}(\lambda) \geq \hat{c}$ при всех $\lambda \in \mathbb{R}_+^m$.

Лемма 3. Пусть есть k аффинных функционалов $l_i(x)$ на \mathbb{R}^n таких, что при $x_0 \in \mathbb{R}_+^n$ для любого $i = 1, \dots, k$ $l_i(x_0) = c_0$.

Пусть функция $l(x) = \max_{1 \leq i \leq k} l_i(x) \geq c_0$ для любого $x \in \mathbb{R}_+^n$. Тогда:

1) существует выпуклая комбинация $\sum_{i=1}^k \hat{\alpha}_i l_i(x) \geq c_0$ при всех $x \in \mathbb{R}_+^n$;

2) отличны от нуля будут не более, чем $n + 1$ координата вектора $\hat{\alpha}$.

Доказательство. Если докажем, что существует какая-то выпуклая комбинация $l_i(x)$, удовлетворяющая утверждению 1, то утверждение 2 следует немедленно по теореме Каратеодори [3].

Введем индикаторную функцию $\delta B(x) = \begin{cases} +\infty, & x \notin B, \\ 0, & x \in B. \end{cases}$ Тогда условие

$l(x) \geq c_0$ при любом $x \in \mathbb{R}_+^n$ эквивалентно неравенству $l(x) + \delta \mathbb{R}_+^n(x) \geq c_0$. Но тогда точка x_0 — точка абсолютного минимума, и по теореме [5, с. 70] $0 \in \partial(l(x_0) + \delta \mathbb{R}_+^n(x_0))$, где символ ∂ обозначает субдифференциал соответствующей функции. По теореме Моро — Рокафеллара [3]

$$0 \in \partial(l(x_0) + \delta \mathbb{R}_+^n(x_0)) = \partial l(x_0) + \partial \delta \mathbb{R}_+^n(x_0).$$

Но тогда по теореме Дубовицкого — Милютинна о субдифференциале максимума [3] легко получаем утверждение 1 леммы. Лемма доказана.

Таким образом, выбранная точка $(1, \lambda^*)$ действительно удовлетворяет свойствам точки $(1, \hat{\lambda})$.

З а м е ч а н и е 3. Из утверждения 2 леммы 3 следует, что достаточно взвешивать не более $m + 1$ стратегий $\pi^*(i)$. Для них, как и для функционалов $\mathcal{L}_1(\pi^*(i), \lambda, 1)$ (см. доказательство теоремы Каратеодори в [3]), справедливо: если $\pi^*(j) = \pi(\alpha^k, \pi^*(1), \dots, \pi^*(k))$, то $\alpha_i = \delta_{ij}$, где δ_{ij} — символ Кронекера.

Итак, доказана следующая теорема.

Т е о р е м а. Пусть множество допустимых стратегий \mathcal{D} непусто, тогда существует оптимальная стратегия π_0 , являющаяся взвесью не более чем $m + 1$ стационарной стратегии $\pi(\hat{\lambda})$, где $\hat{\lambda} \in \mathbb{R}_+$.

1. Фрид Е. Б. Об оптимальных стратегиях в задачах управления с ограничениями // Теория вероятностей и ее применения. — 1972, 17, № 1. — С. 194—199.
2. Неве Ж. Математические основы теории вероятностей. — М.: Мир, 1969. — 310 с.
3. Блекуэлл Д. Динамическое программирование в задачах с затухающим действием // Математика: сб. пер. — 1967. — 11, № 4. — С. 151—160.
4. Алексеев В. М., Тихомиров В. М., Фомин С. В. Оптимальное управление. — М.: Наука, 1979. — 429 с.
5. Алексеев В. М., Галеев Э. М., Тихомиров В. М. Сборник задач по оптимизации. — М.: Наука, 1984. — 288 с.