

# 关于Par VI的若干问题

WangFangyikang

2023 年 3 月 28 日

## 目录

<b>1 Computational Complexity</b>	<b>2</b>
1.1 Value Iteration . . . . .	2
1.2 Value Iteration . . . . .	3
1.3 Linear Programming Approach . . . . .	4
1.3.1 原始问题 . . . . .	4
1.3.2 对偶问题 . . . . .	4
<b>2 Computational Complexity</b>	<b>6</b>
2.1 Value Iteration . . . . .	6
2.2 Value Iteration . . . . .	7
2.3 Linear Programming Approach . . . . .	8
2.3.1 原始问题 . . . . .	8
2.3.2 对偶问题 . . . . .	8

# 1 Computational Complexity

我们定义 $L, (P, r, \gamma)$ 是确定一个MDP所需的bit-size, 假定算术运算 $+$ ,  $-$ ,  $\times$ ,  $\div$ 都使用单位时间(unit time)。我们希望找到一个算法, 它能找到在 $L, (P, r, \gamma)$ 、 $|S|$ 、 $|A|$ 的多项式时间内找到最优策略。如果一个算法能在 $|S|$ 、 $|A|$ 的多项式时间 (与 $L, (P, r, \gamma)$ 无关) 内找到最优策略, 则称该算法为strongly polynomial。

## 1.1 Value Iteration

**引理 1.1** (收缩性). 对任意两个向量 $Q, Q' \in \mathbb{R}^{|S| \times |A|}$ ,  $\|\tau Q - \tau Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$ 。这说明了 $\tau$ 是一个收缩映射。

证明. 首先证明对于所有的 $s \in S$ ,  $|V_Q(s) - V_{Q'}(s)| \leq \max_{a \in A} |Q(s, a) - Q'(s, a)|$ 。

不妨设 $V_Q(s) > V_{Q'}(s)$  (另一半是对称的), 记 $a^* = \arg \max_a Q(s, a)$ 。

$$|V_Q(s) - V_{Q'}(s)| = \max_{a \in A} Q(s, a) - \max_{a \in A} Q'(s, a) \leq Q(s, a^*) - Q'(s, a^*) \leq \max_{a \in A} |Q(s, a) - Q'(s, a)|$$

$$\begin{aligned} \|\tau Q - \tau Q'\|_\infty &= \|r + \gamma P V_Q - r - \gamma P V_{Q'}\|_\infty \\ &= \gamma \|P(V_Q - V_{Q'})\|_\infty \\ &\leq \gamma \|V_Q - V_{Q'}\|_\infty \\ &= \gamma \max_s |V_Q(s) - V_{Q'}(s)| \\ &\leq \gamma \max_s \max_a |Q(s, a) - Q'(s, a)| \\ &= \gamma \|Q - Q'\|_\infty \end{aligned} \tag{1}$$

□

由该引理可见, 经历了 $k$ 次迭代后,  $\|Q_k - Q^*\|_\infty \leq \gamma^k \|Q_0 - Q^*\|_\infty$ , 因此  $\lim_{k \rightarrow \infty} Q_k - Q^* = 0$ , 算法收敛。

**引理 1.2** (Q-Error Amplification). 对任意向量 $Q \in \mathbb{R}^{|S| \times |A|}$ ,  $V^{\tau Q} \geq V^* - \frac{2\|Q - Q^*\|_\infty}{1 - \gamma} \mathbb{1}$

其中 $\mathbb{1}$ 为全1向量。

证明. 对于固定的 $s$ 以及 $a = \pi_Q(s)$ , ( $\pi_Q(s) = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q(s, a)$ )

$$\begin{aligned}
V^*(s) - V^{\pi_Q}(s) &= Q^*(s, \pi^*(s)) - Q^{\pi_Q}(s, a) \\
&= Q^*(s, \pi^*(s)) - Q^*(s, a) + Q^*(s, a) - Q^{\pi_Q}(s, a) \\
&= Q^*(s, \pi^*(s)) - Q^*(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[V^*(s') - V^{\pi_Q}(s')] \\
&\leq Q^*(s, \pi^*(s)) - Q(s, \pi^*(s)) + Q(s, a) - Q^*(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[V^*(s') - V^{\pi_Q}(s')] \\
&\leq 2\|Q - Q^*\|_\infty + \gamma\|V^* - V^{\pi_Q}\|_\infty
\end{aligned} \tag{2}$$

□

**定理 1.1** (Q-value iteration convergence). 设 $Q^{(0)} = 0$ ,  $Q^{(k+1)} = \tau Q^{(k)}$ ,  $k = 0, 1, \dots$ ,  $\pi^{(k) = \pi_{Q^{(k)}}}$ , 当 $k \geq \frac{\log \frac{2}{(1-\gamma)^2 \epsilon}}{1-\gamma}$ ,  $V^{\pi^{(k)}} \geq V^* - \epsilon \mathbb{1}$ , 即 $k$ 次迭代后 $V^{\pi^{(k)}}$ 和 $V^*$ 非常接近。

证明.  $\|Q^{(k)} - Q^*\|_\infty = \|\tau^k Q^{(0)} - \tau^k Q^*\|_\infty \leq \gamma^k \|Q^{(0)} - Q^*\|_\infty = (1 - \gamma)^k \leq \frac{\exp(-(1-\gamma)k)}{1-\gamma}$  □

## 1.2 Value Iteration

策略迭代算法从任意一个策略 $\pi_0$ 出发, 并对 $k=0,1,2,\dots$ 重复接下来的两步: 1.策略评估. 计算 $Q^{\pi_k}$

2.策略提升. 更新策略:  $\pi_{k+1} = \pi_{Q^{\pi_k}}$ , 即 $\pi_{k+1}(s) = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q^{\pi_k}(s, a)$

**引理 1.3.** 1.  $Q^{\pi_{k+1}} \geq \tau Q^{\pi_k} \geq Q^{\pi_k}$

2.  $\|Q^{\pi_{k+1}} - Q^*\|_\infty \leq \gamma \|Q^{\pi_k} - Q^*\|_\infty$

证明. 首先证明 $\tau Q^{\pi_k} \geq Q^{\pi_k}$ . 注意到策略迭代中的策略都是确定策略, 对于所有 $k, s$ ,  $V^{\pi_k}(s) = Q^{\pi_k}(s, \pi_k(s))$ .

$$\begin{aligned}
\tau Q^{\pi_k}(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[\max_{a'} Q^{\pi_k}(s', a')] \\
&\geq r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[Q^{\pi_k}(s', \pi_k(s'))] = Q^{\pi_k}(s, a)
\end{aligned} \tag{3}$$

再证明 $Q^{\pi_{k+1}} \geq \tau Q^{\pi_k}$ , 这需要先证明 $Q^{\pi_{k+1}} \geq Q^{\pi_k}$ :

$$Q^{\pi_k} = r + \gamma P^{\pi_k} Q^{\pi_k} \leq r + \gamma P^{\pi_{k+1}} Q^{\pi_k} \leq \sum_{t=0}^{\infty} \gamma^t (P^{\pi_{k+1}})^t r = Q^{\pi_{k+1}}$$

第一个不等式是因为 $\pi_{k+1}$ 是greedy policy, 所以一定比 $\pi_k$ 更好, 第二个不等式由递归得到. 因此

$$\begin{aligned}
Q^{\pi_{k+1}}(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[Q^{\pi_{k+1}}(s', \pi_{k+1}(s'))] \\
&\geq r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[Q^{\pi_k}(s', \pi_{k+1}(s'))] \\
&= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[\max_{a'} Q^{\pi_k}(s', a')] = \tau Q^{\pi_k}(s, a)
\end{aligned} \tag{4}$$

(1) 式证明完成。下面证明 (2) 式:

$$\|Q^* - Q^{\pi_{k+1}}\|_\infty \leq \|Q^* - \tau Q^{\pi_k}\|_\infty = \|\tau Q^* - \tau Q^{\pi_k}\|_\infty \leq \gamma \|Q^* - Q^{\pi_k}\|_\infty$$

证明完成。  $\square$

**定理 1.2** (Policy iteration convergence). 设  $Q^{\pi_0} = 0, \pi_0$  为任意初始策略。当  $k \geq \frac{\log \frac{1}{(1-\gamma)\epsilon}}{1-\gamma}$ , 第  $k$  个策略有这样的 *performance bound*:

$$Q^{\pi_k} \geq Q^* - \epsilon$$

证明.

$$\|Q^* - Q^{\pi_k}\|_\infty \leq \gamma \|Q^* - Q^{\pi_{k-1}}\|_\infty \leq \gamma^k \|Q^* - Q^{\pi_0}\|_\infty = \gamma^k \|Q^*\|_\infty = (1 - (1-\gamma))^k \|Q^*\|_\infty \leq \frac{\exp(-(1-\gamma)k)}{1-\gamma} \|Q^*\|_\infty \leq \frac{\exp(-(1-\gamma)k)}{1-\gamma} \epsilon$$

(5)  $\square$

### 1.3 Linear Programming Approach

线性规划的方法可以在严格多项式时间内解决问题。

#### 1.3.1 原始问题

最初的想法是求解

$$\begin{aligned} \min_{V \in \mathbb{R}^{|S|}} \quad & \sum_s \mu(s) V(s) \\ \text{subject to } & V(s) \geq \max_{a \in \mathcal{A}} [r(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s')] \quad \forall s \in \mathcal{S} \end{aligned}$$

但这不是LP问题。将其转化为LP问题, 得到:

$$\begin{aligned} \min_{V \in \mathbb{R}^{|S|}} \quad & \sum_s \mu(s) V(s) \\ \text{subject to } & V(s) \geq r(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s') \quad \forall a \in \mathcal{A}, s \in \mathcal{S} \end{aligned}$$

其中,  $\mu(s)$  是初始状态分布。如果  $\mu$  has full support, 那么该问题的唯一解就是  $V^*(s)$ 。

证明. 证明利用了  $\tau$  的单调性, 即  $V \geq V'$  时, 有  $\tau V \geq \tau V'$ 。

令  $V' = \tau V$ , 则  $\tau V \geq \tau V' = \tau^2 V$ 。迭代得到:  $V \geq \tau^\infty = V^*$ 。

因此该约束条件下得到的解都是  $V \geq V^*$  的情况, 由于目标函数是求

$$\min_{V \in \mathbb{R}^{|S|}} \sum_s \mu(s) V(s),$$

最终得到的解即  $V = V^*$ 。  $\square$

#### 1.3.2 对偶问题

对每一个LP都存在一个对偶问题, 原问题的决策变量对应于对偶问题的约束条件, 原问题的约束条件对应于对偶问题的决策变量。

对于固定的策略 $\pi$ ，定义关于状态和动作的visitation measure：

$$d_{s_0}^\pi(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t Pr^\pi(s_t = s, a_t = a | s_0) \quad (6)$$

其中 $Pr^\pi(s_t = s, a_t = a | s_0)$ 是从状态 $s_0$ 出发，经过策略 $\pi$ ，到达 $s_t = s, a_t = a$ 的概率。并记 $d_\mu^\pi(s, a) = \mathbb{E}_{s_0 \sim \mu} [d_{s_0}^\pi(s, a)]$

对于任意的状态 $s$ 有：

$$\sum_a d_\mu^\pi(s, a) = (1 - \gamma)\mu(s) + \gamma \sum_{s', a'} P(s | s', a') d_\mu^\pi(s', a') \quad (7)$$

证明. 左边：

$$\begin{aligned} \sum_a d_\mu^\pi(s, a) &= d_\mu^\pi(s) \\ &= \mathbb{E}_{s_0 \sim \mu} [d_{s_0}^\pi(s)] \\ &= \mathbb{E}_{s_0 \sim \mu} [(1 - \gamma) \sum_{t=0}^{\infty} \gamma^t Pr^\pi(s_t = s | s_0)] \end{aligned} \quad (8)$$

右边：

$$\begin{aligned} &\gamma \sum_{s', a'} P(s | s', a') d_\mu^\pi(s', a') + (1 - \gamma)\mu(s) \\ &= \gamma \sum_{s', a'} P(s | s', a') \mathbb{E}_{s_0 \sim \mu} [(1 - \gamma) \sum_{t=0}^{\infty} \gamma^t Pr^\pi(s_t = s', a_t = a' | s_0)] + (1 - \gamma)\mu(s) \\ &= \mathbb{E}_{s_0 \sim \mu} [(1 - \gamma) \sum_{t=0}^{\infty} \sum_{s', a'} \gamma^{t+1} P(s | s', a') Pr^\pi(s_t = s', a_t = a' | s_0)] + (1 - \gamma)\mu(s) \quad (9) \\ &= \mathbb{E}_{s_0 \sim \mu} [(1 - \gamma) \sum_{t=0}^{\infty} \sum_{s', a'} \gamma^{t+1} Pr^\pi(s | s_0)] + (1 - \gamma)\mu(s) \\ &= \mathbb{E}_{s_0 \sim \mu} [(1 - \gamma) \sum_{t=0}^{\infty} \gamma^t Pr^\pi(s_t = s | s_0)] \end{aligned}$$

所以左边=右边。  $\square$

定义一个状态-动作多面体：

$$\mathcal{K}_\mu := \{d | d \geq 0 \text{ and } \sum_a d(s, a) = (1 - \gamma)\mu(s) + \gamma \sum_{s', a'} P(s | s', a') d(s', a')\}$$

可以看到， $\mathcal{K}_\mu$ 是所有可以的状态-动作分布的集合，即 $d \in \mathcal{K}_\mu$ 当且仅当存在一个平稳的策略 $\pi$ 使得 $d_\mu^\pi = d$ 。

有了这些定义，对偶问题可以写成：

$$\max \quad \frac{1}{1 - \gamma} \sum_{s, a} d_\mu(s, a) r(s, a)$$

subject to  $d \in \mathcal{K}_\mu$

目标函数可以看作是每个状态-动作对的密度乘上奖励，加权求和来求总奖励，求这个总奖励的最大值。解这个对偶LP求得一个解 $d^*$ ，那么就可以求得最优策略：

$$\pi^*(s, a) = \frac{d^*(s, a)}{\sum_{a'} d^*(s, a')}$$

另一种求最优策略的方法是求 $\underset{a}{\operatorname{argmax}} d^*(s, a')$ 。

## 2 Computational Complexity

我们定义 $L, (P, r, \gamma)$ 是确定一个MDP所需的bit-size，假定算术运算 $+$ ， $-$ ， $\times$ ， $\div$ 都使用单位时间(unit time)。我们希望找到一个算法，它能找到在 $L, (P, r, \gamma)$ 、 $|\mathcal{S}|$ 、 $|\mathcal{A}|$ 的多项式时间内找到最优策略。如果一个算法能在 $|\mathcal{S}|$ 、 $|\mathcal{A}|$ 的多项式时间（与 $L, (P, r, \gamma)$ 无关）内找到最优策略，则称该算法为strongly polynomial。

### 2.1 Value Iteration

**引理 2.1** (收缩性). 对任意两个向量 $Q, Q' \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ ,  $\|\tau Q - \tau Q'\|_\infty \leq \gamma \|Q - Q'\|_\infty$ 。这说明了 $\tau$ 是一个收缩映射。

证明. 首先证明对于所有的 $s \in \mathcal{S}$ ,  $|V_Q(s) - V_{Q'}(s)| \leq \max_{a \in \mathcal{A}} |Q(s, a) - Q'(s, a)|$ 。

不妨设 $V_Q(s) > V_{Q'}(s)$  (另一半是对称的)，记 $a^* = \underset{a}{\operatorname{argmax}} Q(s, a)$ 。

$$|V_Q(s) - V_{Q'}(s)| = \max_{a \in \mathcal{A}} Q(s, a) - \max_{a \in \mathcal{A}} Q'(s, a) \leq Q(s, a^*) - Q'(s, a^*) \leq \max_{a \in \mathcal{A}} |Q(s, a) - Q'(s, a)|$$

$$\begin{aligned} \|\tau Q - \tau Q'\|_\infty &= \|r + \gamma P V_Q - r - \gamma P V_{Q'}\|_\infty \\ &= \gamma \|P(V_Q - V_{Q'})\|_\infty \\ &\leq \gamma \|V_Q - V_{Q'}\|_\infty \\ &= \gamma \max_s |V_Q(s) - V_{Q'}(s)| \\ &\leq \gamma \max_s \max_a |Q(s, a) - Q'(s, a)| \\ &= \gamma \|Q - Q'\|_\infty \end{aligned} \tag{10}$$

□

由该引理可见，经历了 $k$ 次迭代后， $\|Q_k - Q^*\|_\infty \leq \gamma^k \|Q_0 - Q^*\|_\infty$ ，因此 $\lim_{k \rightarrow \infty} Q_k - Q^* = 0$ ，算法收敛。

**引理 2.2** (Q-Error Amplification). 对任意向量 $Q \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ ,  $V^{\pi_Q} \geq V^* - \frac{2\|Q - Q^*\|_\infty}{1 - \gamma} \mathbb{1}$

其中 $\mathbb{1}$ 为全1向量。

证明. 对于固定的 $s$ 以及 $a = \pi_Q(s)$ , ( $\pi_Q(s) = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q(s, a)$ )

$$\begin{aligned}
V^*(s) - V^{\pi_Q}(s) &= Q^*(s, \pi^*(s)) - Q^{\pi_Q}(s, a) \\
&= Q^*(s, \pi^*(s)) - Q^*(s, a) + Q^*(s, a) - Q^{\pi_Q}(s, a) \\
&= Q^*(s, \pi^*(s)) - Q^*(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[V^*(s') - V^{\pi_Q}(s')] \\
&\leq Q^*(s, \pi^*(s)) - Q(s, \pi^*(s)) + Q(s, a) - Q^*(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[V^*(s') - V^{\pi_Q}(s')] \\
&\leq 2\|Q - Q^*\|_\infty + \gamma\|V^* - V^{\pi_Q}\|_\infty
\end{aligned} \tag{11}$$

□

**定理 2.1** (Q-value iteration convergence). 设 $Q^{(0)} = 0$ ,  $Q^{(k+1)} = \tau Q^{(k)}$ ,  $k = 0, 1, \dots$ ,  $\pi^{(k) = \pi_{Q^{(k)}}}$ , 当 $k \geq \frac{\log \frac{2}{(1-\gamma)^2 \epsilon}}{1-\gamma}$ ,  $V^{\pi^{(k)}} \geq V^* - \epsilon \mathbb{1}$ , 即 $k$ 次迭代后 $V^{\pi^{(k)}}$ 和 $V^*$ 非常接近。

证明.  $\|Q^{(k)} - Q^*\|_\infty = \|\tau^k Q^{(0)} - \tau^k Q^*\|_\infty \leq \gamma^k \|Q^{(0)} - Q^*\|_\infty = (1 - \gamma)^k \leq \frac{\exp(-(1-\gamma)k)}{1-\gamma}$  □

## 2.2 Value Iteration

策略迭代算法从任意一个策略 $\pi_0$ 出发, 并对 $k=0,1,2,\dots$ 重复接下来的两步: 1.策略评估. 计算 $Q^{\pi_k}$

2.策略提升. 更新策略:  $\pi_{k+1} = \pi_{Q^{\pi_k}}$ , 即 $\pi_{k+1}(s) = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q^{\pi_k}(s, a)$

**引理 2.3.** 1.  $Q^{\pi_{k+1}} \geq \tau Q^{\pi_k} \geq Q^{\pi_k}$

2.  $\|Q^{\pi_{k+1}} - Q^*\|_\infty \leq \gamma \|Q^{\pi_k} - Q^*\|_\infty$

证明. 首先证明 $\tau Q^{\pi_k} \geq Q^{\pi_k}$ 。注意到策略迭代中的策略都是确定策略, 对于所有 $k, s$ ,  $V^{\pi_k}(s) = Q^{\pi_k}(s, \pi_k(s))$ 。

$$\begin{aligned}
\tau Q^{\pi_k}(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[\max_{a'} Q^{\pi_k}(s', a')] \\
&\geq r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[Q^{\pi_k}(s', \pi_k(s'))] = Q^{\pi_k}(s, a)
\end{aligned} \tag{12}$$

再证明 $Q^{\pi_{k+1}} \geq \tau Q^{\pi_k}$ , 这需要先证明 $Q^{\pi_{k+1}} \geq Q^{\pi_k}$ :

$$Q^{\pi_k} = r + \gamma P^{\pi_k} Q^{\pi_k} \leq r + \gamma P^{\pi_{k+1}} Q^{\pi_k} \leq \sum_{t=0}^{\infty} \gamma^t (P^{\pi_{k+1}})^t r = Q^{\pi_{k+1}}$$

第一个不等式是因为 $\pi_{k+1}$ 是greedy policy, 所以一定比 $\pi_k$ 更好, 第二个不等式由递归得到。因此

$$\begin{aligned}
Q^{\pi_{k+1}}(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[Q^{\pi_{k+1}}(s', \pi_{k+1}(s'))] \\
&\geq r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[Q^{\pi_k}(s', \pi_{k+1}(s'))] \\
&= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[\max_{a'} Q^{\pi_k}(s', a')] = \tau Q^{\pi_k}(s, a)
\end{aligned} \tag{13}$$

(1) 式证明完成。下面证明 (2) 式:

$$\|Q^* - Q^{\pi_{k+1}}\|_\infty \leq \|Q^* - \tau Q^{\pi_k}\|_\infty = \|\tau Q^* - \tau Q^{\pi_k}\|_\infty \leq \gamma \|Q^* - Q^{\pi_k}\|_\infty$$

证明完成。  $\square$

**定理 2.2** (Policy iteration convergence). 设  $Q^{\pi_0} = 0, \pi_0$  为任意初始策略。当  $k \geq \frac{\log \frac{1}{(1-\gamma)\epsilon}}{1-\gamma}$ , 第  $k$  个策略有这样的 *performance bound*:

$$Q^{\pi_k} \geq Q^* - \epsilon$$

证明.

$$\|Q^* - Q^{\pi_k}\|_\infty \leq \gamma \|Q^* - Q^{\pi_{k-1}}\|_\infty \leq \gamma^k \|Q^* - Q^{\pi_0}\|_\infty = \gamma^k \|Q^*\|_\infty = (1 - (1-\gamma))^k \|Q^*\|_\infty \leq \frac{\exp(-(1-\gamma)k)}{1-\gamma} \|Q^*\|_\infty \leq \frac{\exp(-(1-\gamma)k)}{1-\gamma} \epsilon$$

(14)  $\square$

## 2.3 Linear Programming Approach

线性规划的方法可以在严格多项式时间内解决问题。

### 2.3.1 原始问题

最初的想法是求解

$$\begin{aligned} \min_{V \in \mathbb{R}^{|S|}} \quad & \sum_s \mu(s) V(s) \\ \text{subject to } & V(s) \geq \max_{a \in \mathcal{A}} [r(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s')] \quad \forall s \in \mathcal{S} \end{aligned}$$

但这不是LP问题。将其转化为LP问题, 得到:

$$\begin{aligned} \min_{V \in \mathbb{R}^{|S|}} \quad & \sum_s \mu(s) V(s) \\ \text{subject to } & V(s) \geq r(s, a) + \gamma \sum_{s'} P(s'|s, a) V(s') \quad \forall a \in \mathcal{A}, s \in \mathcal{S} \end{aligned}$$

其中,  $\mu(s)$  是初始状态分布。如果  $\mu$  has full support, 那么该问题的唯一解就是  $V^*(s)$ 。

证明. 证明利用了  $\tau$  的单调性, 即  $V \geq V'$  时, 有  $\tau V \geq \tau V'$ 。

令  $V' = \tau V$ , 则  $\tau V \geq \tau V' = \tau^2 V$ 。迭代得到:  $V \geq \tau^\infty = V^*$ 。

因此该约束条件下得到的解都是  $V \geq V^*$  的情况, 由于目标函数是求

$$\min_{V \in \mathbb{R}^{|S|}} \sum_s \mu(s) V(s),$$

最终得到的解即  $V = V^*$ 。  $\square$

### 2.3.2 对偶问题

对每一个LP都存在一个对偶问题, 原问题的决策变量对应于对偶问题的约束条件, 原问题的约束条件对应于对偶问题的决策变量。



对于固定的策略 $\pi$ ，定义关于状态和动作的visitation measure:

$$d_{s_0}^\pi(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t Pr^\pi(s_t = s, a_t = a | s_0) \quad (15)$$

其中 $Pr^\pi(s_t = s, a_t = a | s_0)$ 是从状态 $s_0$ 出发，经过策略 $\pi$ ，到达 $s_t = s, a_t = a$ 的概率。并记 $d_\mu^\pi(s, a) = \mathbb{E}_{s_0 \sim \mu} [d_{s_0}^\pi(s, a)]$

对于任意的状态 $s$ 有:

$$\sum_a d_\mu^\pi(s, a) = (1 - \gamma)\mu(s) + \gamma \sum_{s', a'} P(s | s', a') d_\mu^\pi(s', a') \quad (16)$$

证明. 左边:

$$\begin{aligned} \sum_a d_\mu^\pi(s, a) &= d_\mu^\pi(s) \\ &= \mathbb{E}_{s_0 \sim \mu} [d_{s_0}^\pi(s)] \\ &= \mathbb{E}_{s_0 \sim \mu} [(1 - \gamma) \sum_{t=0}^{\infty} \gamma^t Pr^\pi(s_t = s | s_0)] \end{aligned} \quad (17)$$

右边:

$$\begin{aligned} &\gamma \sum_{s', a'} P(s | s', a') d_\mu^\pi(s', a') + (1 - \gamma)\mu(s) \\ &= \gamma \sum_{s', a'} P(s | s', a') \mathbb{E}_{s_0 \sim \mu} [(1 - \gamma) \sum_{t=0}^{\infty} \gamma^t Pr^\pi(s_t = s', a_t = a' | s_0)] + (1 - \gamma)\mu(s) \\ &= \mathbb{E}_{s_0 \sim \mu} [(1 - \gamma) \sum_{t=0}^{\infty} \sum_{s', a'} \gamma^{t+1} P(s | s', a') Pr^\pi(s_t = s', a_t = a' | s_0)] + (1 - \gamma)\mu(s) \quad (18) \\ &= \mathbb{E}_{s_0 \sim \mu} [(1 - \gamma) \sum_{t=0}^{\infty} \sum_{s', a'} \gamma^{t+1} Pr^\pi(s | s_0)] + (1 - \gamma)\mu(s) \\ &= \mathbb{E}_{s_0 \sim \mu} [(1 - \gamma) \sum_{t=0}^{\infty} \gamma^t Pr^\pi(s_t = s | s_0)] \end{aligned}$$

所以左边=右边。  $\square$

定义一个状态-动作多面体:

$$\mathcal{K}_\mu := \{d | d \geq 0 \text{ and } \sum_a d(s, a) = (1 - \gamma)\mu(s) + \gamma \sum_{s', a'} P(s | s', a') d(s', a')\}$$

可以看到， $\mathcal{K}_\mu$ 是所有可以的状态-动作分别的集合，即 $d \in \mathcal{K}_\mu$ 当且仅当存在一个平稳的策略 $\pi$ 使得 $d_\mu^\pi = d$ 。

有了这些定义，对偶问题可以写成:

$$\max \quad \frac{1}{1 - \gamma} \sum_{s, a} d_\mu(s, a) r(s, a)$$

subject to  $d \in \mathcal{K}_\mu$

目标函数可以看作是每个状态-动作对的密度乘上奖励，加权求和来求总奖励，求这个总奖励的最大值。解这个对偶LP求得一个解 $d^*$ ,那么就可以求得最优策略：

$$\pi^*(s, a) = \frac{d^*(s, a)}{\sum_{a'} d^*(s, a')}$$

另一种求最优策略的方法是求 $\underset{a}{\operatorname{argmax}} \quad d^*(s, a')$ 。