RL: Доказательство существования и единственности решения уравнений Беллмана

Утверждение

□ В конечном Марковском процессе принятия решений (MDP) оптимальная

политика определяется как политика, которая максимизирует ценность

всех состояний одновременно

Уравнения Беллмана

$$V_{\pi}(s_t) = E_{\pi}(r_t + \gamma V(s_{t+1}))$$

$$Q_{\pi}(s_t) = E_{\pi}(r_t + \gamma Q(s_{t+1}, a_{t+1}))$$

□ Связь между функцией ценности и ценности действия:

$$V_{\pi}(s_t) = E_{\pi}(Q(s_t, a_t))$$

Оптимальная политика

□ Политика является оптимальной тогда и только тогда, когда:

$$V_{\pi^*}(s) \ge V_{\pi}(s) \tag{1}$$

Для любого состояния s и любой другой политики

Теорема

□ Если существует $s^* \in S$ и $a^* \in A(s^*)$ такие, что $Q_{\pi}(s^*, a^*) \ge V_{\pi}(s^*)$, тогда для политики

 $V_{\pi'}(s) \ge V_{\pi}(s) \quad \forall s \in S$

$$\pi'(a|s) = \pi(a|s) (1 - \delta_{s,s^*}) + \delta_{a,a^*} \delta_{s,s^*}$$
 мы имеем

Доказательство: для любого $s \in S$

$$\begin{split} V_{\pi}(s) &= E_{\pi}(Q_{\pi}(s,A_{t})|S_{t} = s) \leq \\ &\leq E_{\pi\prime}(Q_{\pi}(s,A_{t})|S_{t} = s) \\ &= E_{\pi\prime}\left(R_{t+1} + \gamma V_{\pi}(S_{t+1})|S_{t} = s\right) \\ &\leq E_{\pi\prime}(R_{t+1} + \gamma R_{t+2} + \gamma^{2} V_{\pi}(S_{t+2})|S_{t} = s) \\ &\leq E_{\pi\prime}(R_{t+1} + \gamma R_{t+2} + \gamma^{2} R_{t+3} + \gamma^{3} R_{t+4} + \cdots |S_{t} = s) \\ &= V_{\pi\prime}\left(s\right) \end{split}$$

□ Теорема утверждает, что всякий раз, когда есть пара состояния и действия (s*, a*) со значением больше, чем значение состояния s*, относительно политики т, то существует другая политика т', которая лучше или равна (в терминах значений состояния) т во всех состояниях. В результате, если существует оптимальная политика π^* , ее значения должны удовлетворять

$$V_{\pi^*}(s) = \max_{a \in A(s)} Q_{\pi^*}(s, a)$$

Уравнения оптимальности Беллмана

$$V_{\pi^*}(s) = \max_{a \in A(s)} Q_{\pi^*}(s, a)$$

$$V_{\pi^*}(s) = \max_{a \in A} E(R_{t+1} + \gamma V_{\pi^*}(S_{t+1}) | S_t = s, A_t = a)$$

$$= \max_{a \in A} \left(\sum_{r} rp(r|s,a) + \gamma \sum_{s'} V_{\pi^*}(s')p(s'|s,a) \right)$$

Существование и единственность

Чтобы показать, что существует оптимальная политика, необходимо доказать следующие два утверждения:

- 1. Система уравнений оптимальности Беллмана имеет решения, и
- 2. Одно из его решений имеет значения больше или равные значениям других решений во всех состояниях.

Оператор оптимальности Беллмана

$$V = \{V(s) : s \in S\} \in R^{|S|}$$

□ Оператор оптимальности Беллмана

$$B: R^{|S|} \to R^{|S|}$$
 такой, что

$$BV(s) = \max_{a \in A(s)} E(R_{t+1} + \gamma V(S_{t+1}) | S_t = s, A_t = a)$$

$$V_{\pi^*} = BV_{\pi^*}$$

Сжимающие отображения и теорема Банах о неподвижной точек

Теорема:

Пусть (X, d) - непустое полное метрическое пространство

Пусть $B: X \to X$ сжимающее отображение на X, т.е. существует число $0 \le \alpha < 1$,

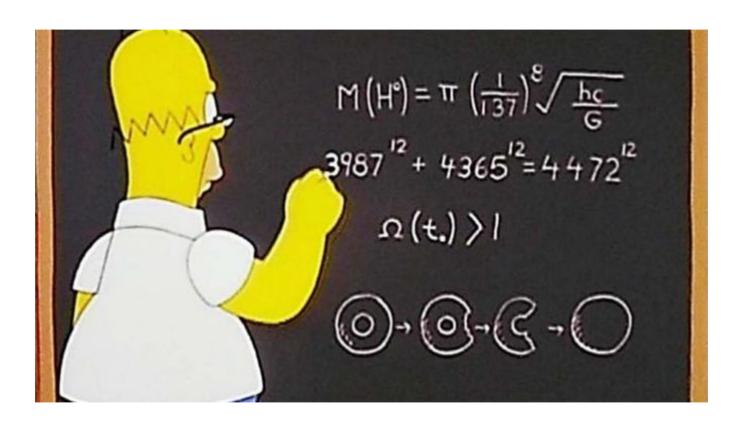
такое что

$$d(Bx, By) \le \alpha d(x, y)$$
, для всех x и y из X

Тогда у отображения B существует, и притом единственная, неподвижная

точка x^* из X (т.е. $Bx^* = x^*$)

Доказательство



Оператор оптимальности Беллмана – это сжимающее отображение в бесконечной норме

Бесконечная норма:

$$||V - V'||_{\infty} = \max_{s \in S} |V(s) - V(s')|$$

Лемма: Пусть задано ограниченное множество A, тогда для любых значений функций f, таких, что $f_1 \colon A \to R$ и $f_2 \colon A \to R$ справедливо:

$$\left[\max_{a \in A} f_1(a) - \max_{a \in A} f_2(a) \right] \le \max_{a \in A} [f_1(a) - f_2(a)]$$

Теорема

На метрическом пространстве $(R^{|S|}, \|\cdot\|_{\infty})$, оператор оптимальности Беллмана

 $B: R^{|S|} \to R^{|S|}$ сжимающий, т.е.

$$\exists \alpha \in [0,1): \|BV - BV'\|_{\infty} \le \alpha \|V - V'\|_{\infty} \ \forall V, V' \in R^{|S|}$$

Доказательство

$$\begin{split} \|BV - BV'\|_{\infty} &= \max_{s \in S} |BV(s) - BV'(s)| \\ &= \max_{s \in S} \left| \max_{a \in A(s)} E[R_{t+1} + \gamma V(s_{t+1})] - \max_{a \in A(s)} E[R_{t+1} + \gamma V'(s_{t+1})] \right| \\ &\leq \max_{s \in S} \max_{a \in A(s)} |E[R_{t+1} + \gamma V(s_{t+1})] - E[R_{t+1} + \gamma V'(s_{t+1})]| \\ &= \max_{s \in S} \max_{a \in A(s)} |\gamma E[V(s_{t+1}) - V'(s_{t+1})]| \\ &\leq \max_{s \in S} \max_{a \in A(s)} \gamma E|[V(s_{t+1}) - V'(s_{t+1})]| \\ &\leq \max_{s \in S} \max_{a \in A(s)} \gamma E|[V(s_{t+1}) - V'(s_{t+1})]| \\ &\leq \max_{s \in S} \max_{a \in A(s)} \gamma E|[V(s_{t+1}) - V'(s_{t+1})]| \end{split}$$

