

Markov Decision Process II

2023年10月5日 星期四 下午9:31

定理: V 是 Bellman 方程至多一的解. 更进一步的, 如果 $J^*(x) \in \operatorname{Argmax}_{a \in A} \{R(x, a) + \sum_{z \in S} P_a(x, z) V(z)\}$

则有 $\alpha^*(n, x) = J^*(x)$ 是一个最优控制.

证明 (接上次): $U(x, \alpha) = E_x \left[\sum_{n=0}^{\infty} \lambda^n R(x_n^\alpha, \alpha(n, x_n^\alpha)) \right]$

$$E_x \sum_{n=0}^{\infty} \lambda^n V(x_n^{\alpha^*}) = E_x \sum_{n=0}^{\infty} \lambda^n \left[R(x_n^{\alpha^*}, \alpha^*(n, x_n^{\alpha^*})) + \lambda \sum_{z \in S} P_{\alpha^*}(x_n^{\alpha^*}, z) V(z) \right].$$

$$= U(x, \alpha^*) + E_x \sum_{n=0}^{\infty} \lambda^{n+1} E_x [V(x_{n+1}^{\alpha^*}) | x_n^{\alpha^*}]$$

$$= U(x, \alpha^*) + \sum_{n=0}^{\infty} \lambda^{n+1} E_x [V(x_{n+1}^{\alpha^*})].$$

$$\Rightarrow V(x) = U(x, \alpha^*).$$



例子: 假设代理人有随机的资金 x_n . 在每个时间单位消耗了一些, 把消耗量视为控制. R 定义了代理人的

消耗资金的功效 代理人希望最大化 $E_x \left[\sum_{n=0}^{\infty} \frac{1}{(1+r)^n} R(x_n^\alpha, \alpha(n, x_n^\alpha)) \right]$, r 是利率.

$$\text{Bellman 方程: } V(x) = \sup_a \left\{ R(x, a) + \sum_{z \in S} p(x, a, z) V(z) \right\}$$