

## 第二课 马尔可夫决策过程 下

### 1. 预测 (*prediction*) 和控制 (*control*)

- 预测: 给定  $MDP(S, A, P, R, \gamma)$  和策略  $\pi$ , 求解价值函数  $v^\pi$
- 控制: 给定  $MDP(S, A, P, R, \gamma)$ , 求最佳的价值函数  $v^*$  和最佳策略  $\pi$

### 2. 动态规划的特点

- 将问题分解为最佳子结构
  - 最优性原理适用
  - 最优解可被分解为子问题 (最优性原理: 多阶段决策过程的最优决策序列具有这样的性质: 不论初始状态和初始决策如何, 对前面决策所造成的某一状态而言, 其后各阶段的决策序列必须构成最优决策)
- 重叠子问题
  - 子问题多次重复出现
  - 解能被缓存和重新使用

### 3. $MDP$ 满足这两个性质:

- *Bellman* 等式给出了一个迭代的分解
- 价值函数存储并能重复利用解

### 4. $MDP$ 的策略评估

- 目标: 评估  $MDP$  的一个给定策略
- 输出: 策略的价值函数  $v^\pi$
- 具体算法: *synchronous backup*

$$V_{t+1}(s) = \sum_{a \in A} \pi(a | s) \left( R(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) v_t(s') \right)$$

或化为  $MRP(S, P^\pi, R, \gamma)$

$$v_{t+1}(s) = R^\pi(s) + \gamma \sum_{s' \in S} P^\pi(s' | s) v_t(s')$$

### 5. 最佳价值函数与最佳策略

- 最佳价值函数为  $v^* = \max_{\pi} v^\pi(s)$
- 最佳策略为  $\pi^*(s) = \operatorname{argmax}_{\pi} v^\pi(s)$

### 6. $MDP$ 控制

- 算法1: *policy iteration* + *policy improvement*
  - *policy iteration*
    - 评估策略  $\pi$
    - 改进策略:  $\pi' = \operatorname{greedy}(v^\pi)$
    - 上述两步不断迭代最终  $\pi'$  会收敛
  - *policy improvement*
    - 计算  $q^{\pi_i}(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) v^{\pi_i}(s')$
    - 计算  $\pi_{i+1}(s) = \operatorname{argmax}_a q^{\pi_i}(s, a)$
    - *Bellman* 最优等式  $v^\pi = \max_{a \in A} q^\pi(s, a)$ , 满足后可有最优价值函数且

$$v^*(s) = \max_a R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) v^*(s')$$

$$q^*(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) \max_{a'} q^*(s', a')$$

○ 算法2: 值迭代

- 将Bellman最优等式作为更新规则

$$v(s) \leftarrow \max_{a \in A} R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) v(s')$$

- 具体的:

$$k = 1, v_0(s) = 0$$

$$\text{while } k \leq H :$$

for s :

$$q_{k+1}(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) v_k(s')$$

$$v_{k+1}(s) = \max_a q_{k+1}(s, a)$$

$$k \leftarrow k + 1$$

- 得到 $v^*$ 的同时可以得到 $\pi(s) = \argmax_a R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) v^*(s')$