

# 彪哥带你学强化学习

## 5、彻底理解贝尔曼方程

DEEPLY UNDERSTAND REINFORCEMENT LEARNING

讲师：韩路彪



## 贝尔曼方程

意义：可以在后面步骤没走完的情况下计算前面步骤的价值，为时序差分、动态规划等算法提供了理论指导

是什么

1. 用后一步的价值表示前一步的价值
2. 类似递推公式
3. 有多种不同表达形式

$$V^{\pi}(s_t) = R_t + \gamma \mathbb{E}_{S_{t+1}}[V^{\pi}(S_{t+1}) | S_t = s_t]$$

比如：

$V_{t+1} \rightarrow Q_t$

$Q_{t+1} \rightarrow V_t$

$V_{t+1} \rightarrow V_t$

$Q_{t+1} \rightarrow Q_t$

$Q^*_{t+1} \rightarrow Q^*_t$

$V^*_{t+1} \rightarrow V^*_t$

为什么学

1. 是后面算法的基础
2. 可以加深对模型的理解

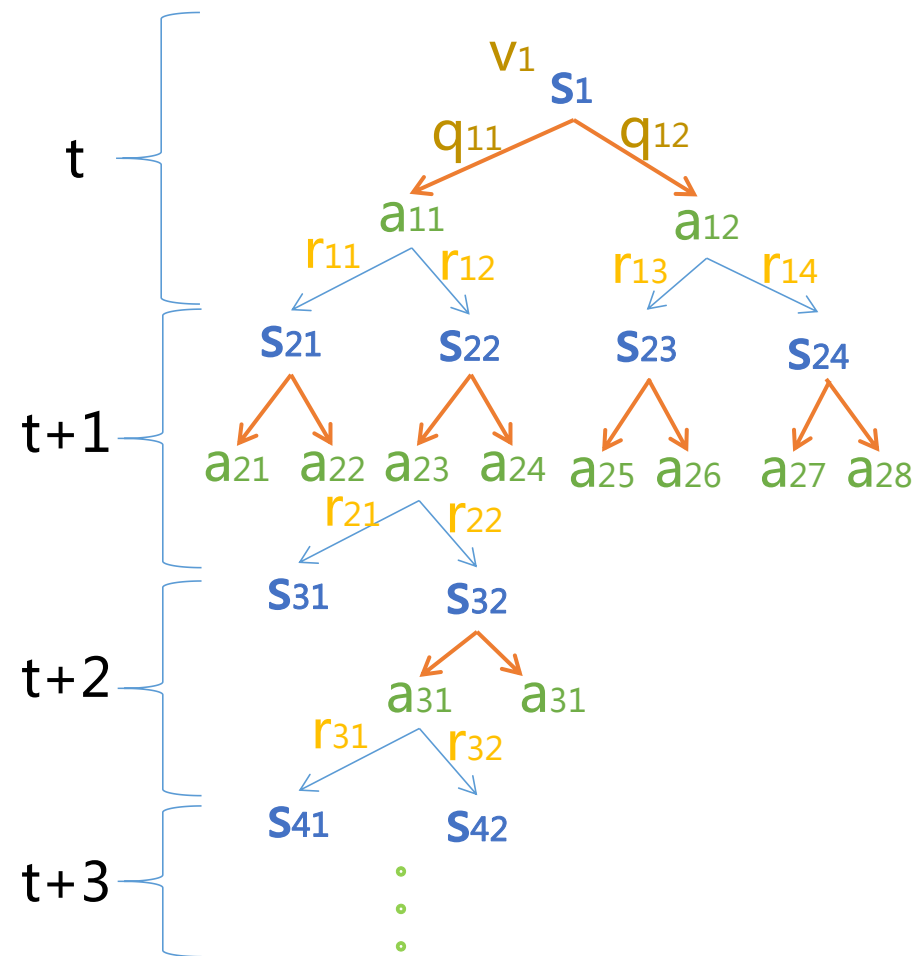
怎么学

以理解含义为主，要会数学公式推导

## 贝尔曼方程 $V_{t+1}$ 到 $Q_t$

$$Q^\pi(s_t, a_t) = \mathbb{E}_{S_{t+1} \sim p(\cdot | s_t, a_t)} [R_t + \gamma V^\pi(S_{t+1})]$$

$$\begin{aligned} Q^\pi(s_t, a_t) &\triangleq \mathbb{E}^\pi(G_t | S_t = s_t, A_t = a_t) \\ &= \mathbb{E}^\pi(R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots | S_t = s_t, A_t = a_t) \\ &= \mathbb{E}^\pi(R_t + \gamma(R_{t+1} + \gamma R_{t+2} + \dots) | S_t = s_t, A_t = a_t) \\ &= \mathbb{E}^\pi(R_t + \gamma V^\pi(S_{t+1}) | S_t = s_t, A_t = a_t) \\ &= \mathbb{E}_{S_{t+1} \sim p(\cdot | s_t, a_t)} [R_t + \gamma V^\pi(S_{t+1})] \end{aligned}$$



贝尔曼方程  $V_{t+1}$  到  $Q_t$

$$Q^\pi(s_t, a_t) = \mathbb{E}_{S_{t+1} \sim p(\cdot | s_t, a_t)} [R_t + \gamma V^\pi(S_{t+1})]$$

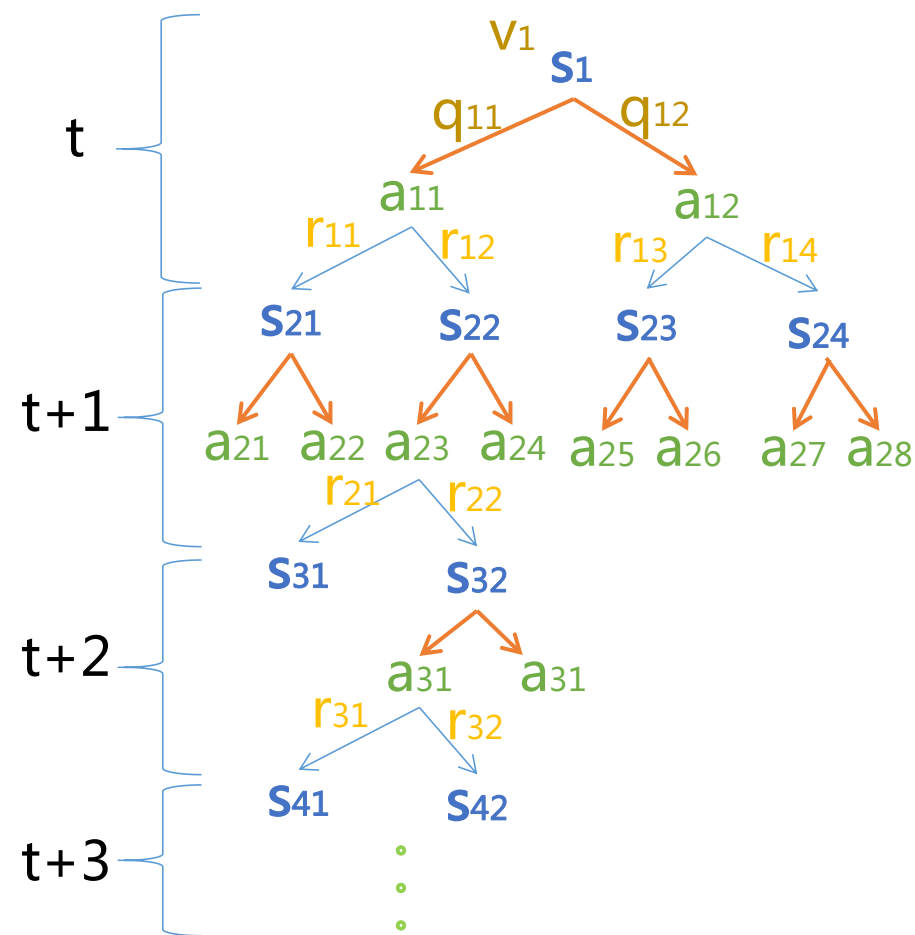
$$Q^\pi(s_t, a_t) = \mathbb{E}_{S_{t+1} \sim p(\cdot | s_t, a_t)} [R(s_t, a_t, S_{t+1}) + \gamma V^\pi(S_{t+1})]$$

$$Q^\pi(s_t, a_t) = \mathbb{E}_{S_{t+1}} [R(s_t, a_t, S_{t+1}) + \gamma V^\pi(S_{t+1}) | S_t = s_t, A_t = a_t]$$

$$Q^\pi(s_t, a_t) = \sum_{S_{t+1} \sim p(\cdot | s_t, a_t)} p(S_{t+1} | s_t, a_t) [R(s_t, a_t, S_{t+1}) + \gamma V^\pi(S_{t+1})]$$

贝尔曼方程  $Q_{t+1}$  到  $V_t$

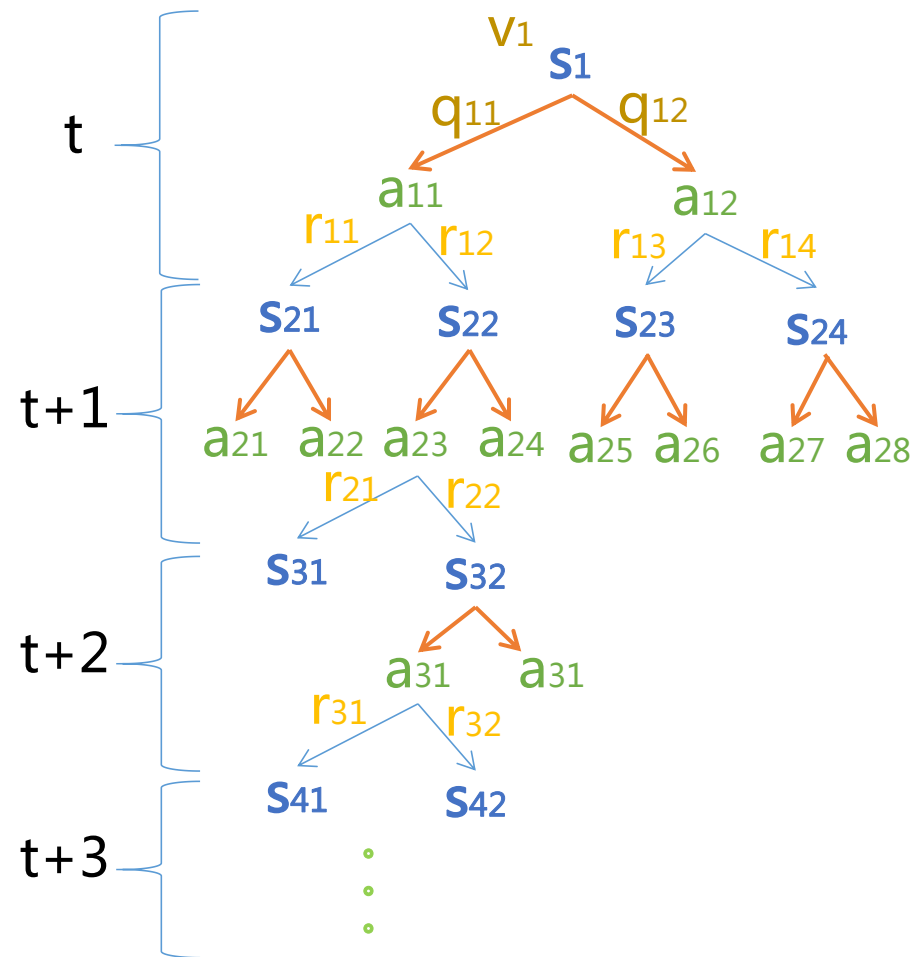
$$V^\pi(s_t) = R_t + \gamma \mathbb{E}_{S_{t+1}, A_{t+1}} [Q^\pi(S_{t+1}, A_{t+1}) | S_t = s_t]$$



## 贝尔曼方程 $V_{t+1}$ 到 $V_t$

$$V^\pi(s_t) = R_t + \gamma \mathbb{E}_{S_{t+1}}[V^\pi(S_{t+1}) | S_t = s_t]$$

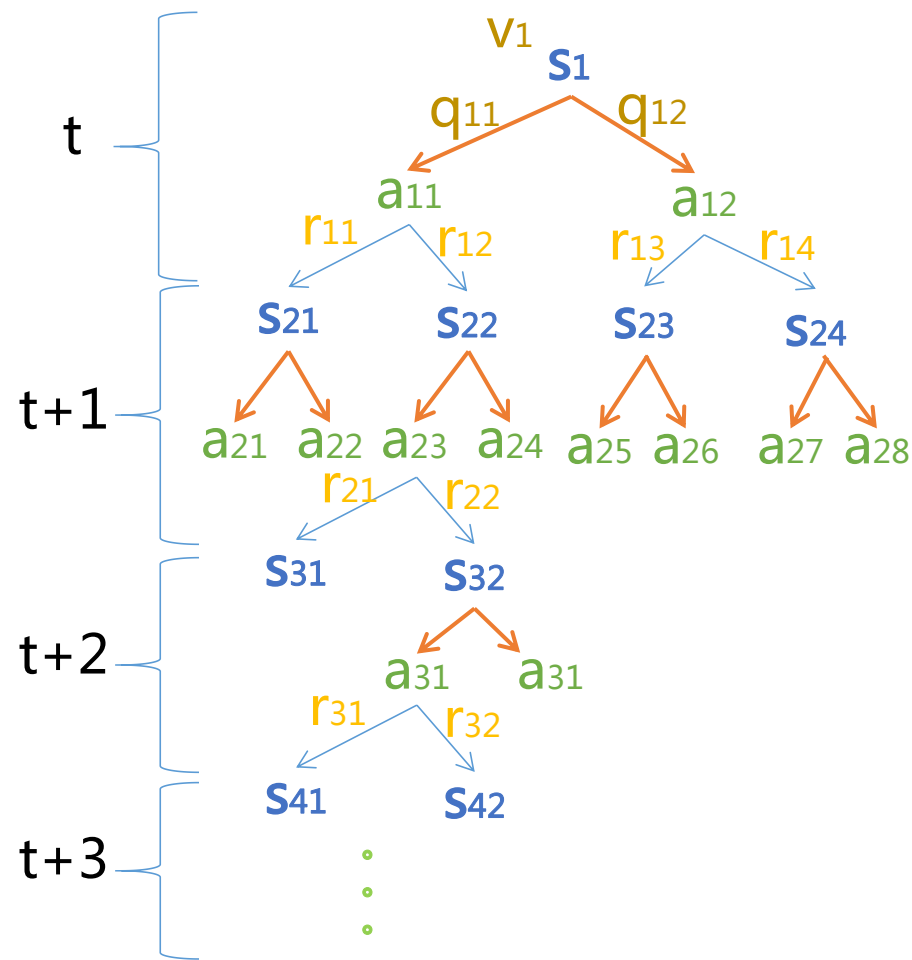
$$V^\pi(s_t) = \mathbb{E}_{A_t, S_{t+1}}[R(s_t, A_t, S_{t+1}) + \gamma V^\pi(S_{t+1}) | S_t = s_t]$$





## 贝尔曼方程 $Q_{t+1}$ 到 $Q_t$

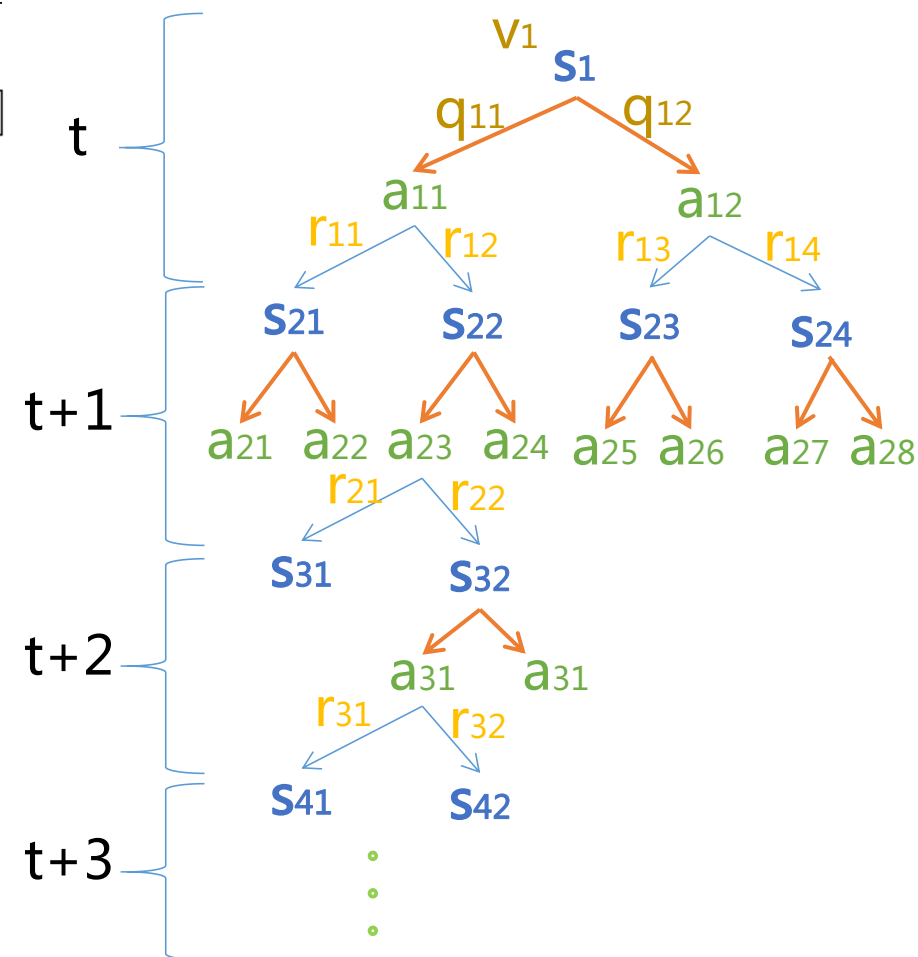
$$Q^\pi(s_t, a_t) = \mathbb{E}_{S_{t+1}, A_{t+1}}^\pi [R(s_t, a_t, S_{t+1}) + \gamma Q^\pi(S_{t+1}, A_{t+1}) | S_t = s_t, A_t = a_t]$$



## 贝尔曼最优方程 $Q^*_{t+1}$ 到 $Q^*_t$

$$Q^*(s_t, a_t) = \mathbb{E}_{S_{t+1} \sim p(\cdot | s_t, a_t)} [R(s_t, a_t, S_{t+1}) + \gamma \max_{A_{t+1} \in \mathcal{A}} Q^*(S_{t+1}, A_{t+1})]$$

$$Q^*(s_t, a_t) = \sum_{S_{t+1} \sim p(\cdot | s_t, a_t)} p(S_{t+1} | s_t, a_t) [R(s_t, a_t, S_{t+1}) + \gamma \max_{A_{t+1} \in \mathcal{A}} Q^*(S_{t+1}, A_{t+1})]$$





贝尔曼最优方程  $V^*_{t+1}$  到  $V^*_t$

$$V^*(s_t) = \max_{A_t \in \mathcal{A}} \left( \sum_{S_{t+1} \sim p(\cdot | s_t, A_t)} p(S_{t+1} | s_t, A_t) [R(s_t, A_t, S_{t+1}) + \gamma V^*(S_{t+1})] \right)$$

