

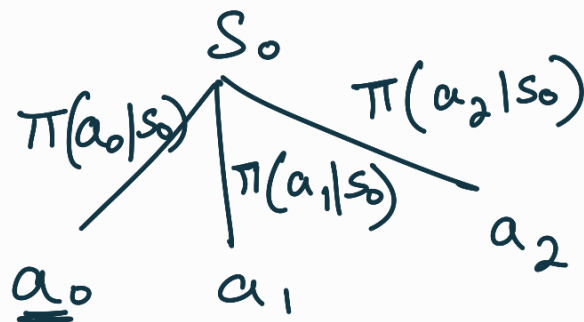
$$\underline{v_\pi(s)} = \underline{E[G_t | S_t = s]}$$

$$\underline{q_\pi(s, a)} = E[G_t | S_t = s, A_t = a]$$

$$S = \{s_0, s_1\}$$

$$A = \{a_0, a_1, a_2\}$$

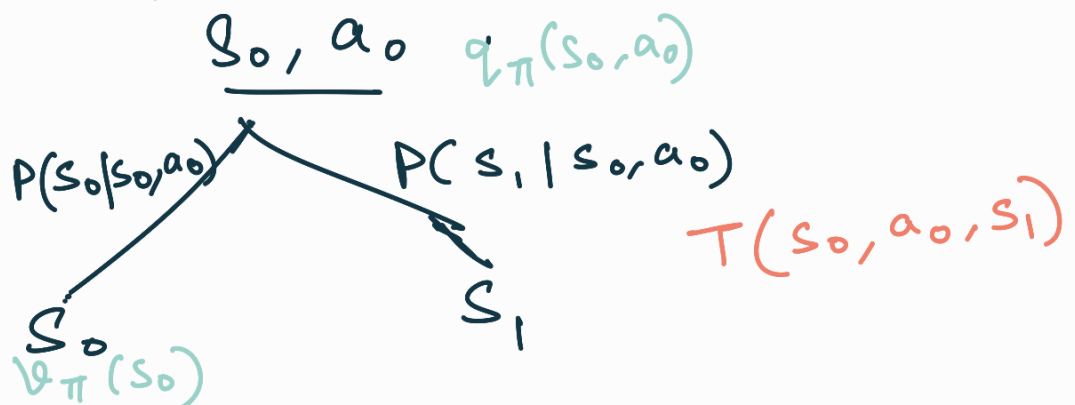
* one level back up diagram starting from state s_0 .



$$\begin{aligned} \underline{v_\pi(s_0)} &= \pi(a_0 | s_0) q_\pi(s_0, a_0) \\ &\quad + \pi(a_1 | s_0) q_\pi(s_0, a_1) \\ &\quad + \pi(a_2 | s_0) q_\pi(s_0, a_2) \end{aligned}$$

$$= \sum_{a \in A} \pi(a | s_0) q_\pi(s_0, a) \rightarrow \textcircled{1}$$

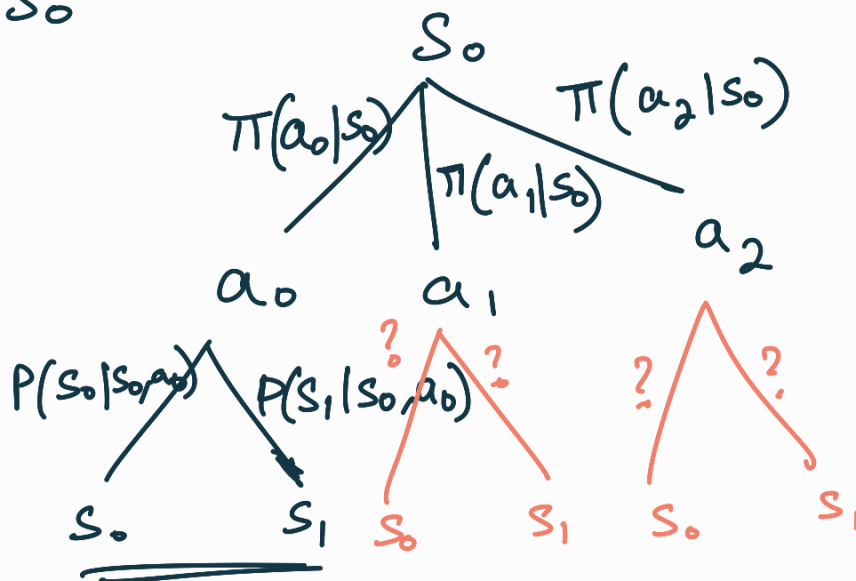
* Starting from state s_0 and taking action a_0 .



$$\begin{aligned}
 q_{\pi}(s_0, a_0) &= P(s_0 | s_0, a_0) \left(r(s_0, a_0, s_0) + \gamma V_{\pi}(s_0) \right) \\
 &\quad + P(s_1 | s_0, a_0) \left(r(s_0, a_0, s_1) + \gamma V_{\pi}(s_1) \right) \\
 &= \sum_{s' \in S} P(s' | s_0, a_0) \left[r(s_0, a_0, s') + \gamma V_{\pi}(s') \right]
 \end{aligned}$$

2-Step backup diagram

* s_0



From (i)

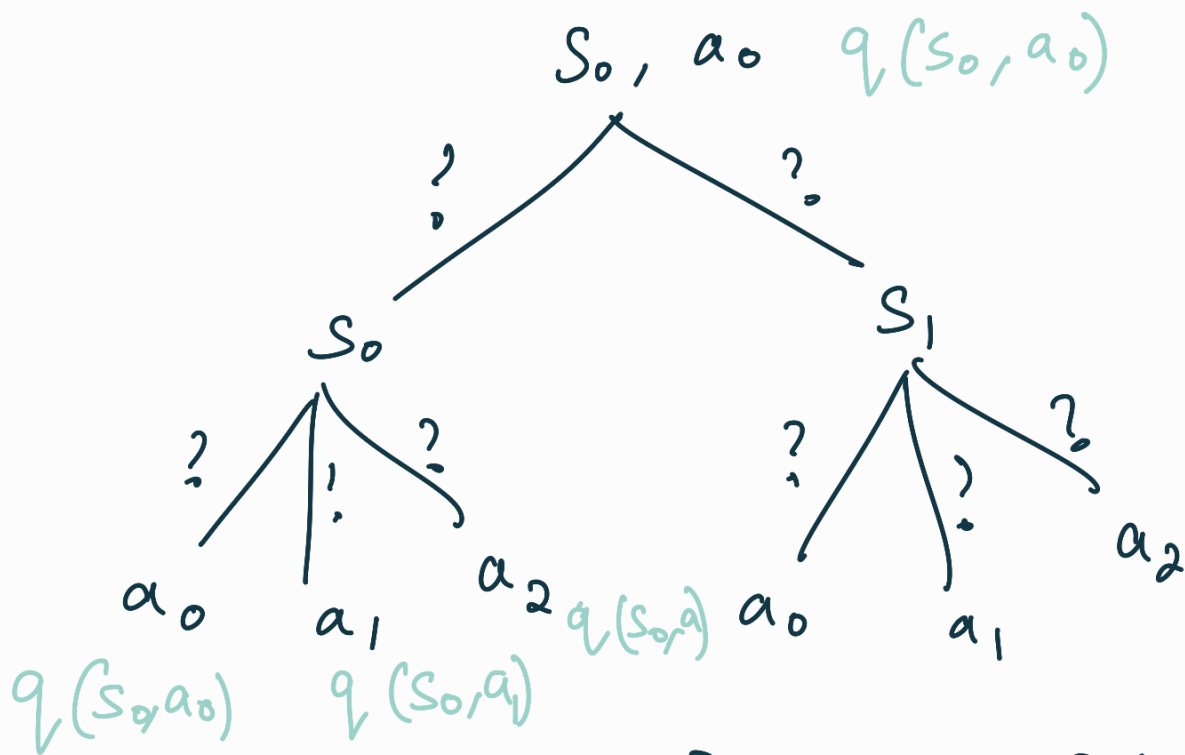
$$V_{\pi}(s_0) = \sum_{a \in A} \pi(a | s_0) q_{\pi}(s_0, a)$$

$$= \sum_{a \in A} \pi(a | s_0) \left[P(s_0 | s_0, a) \left(r(s_0, a, s_0) + \gamma V_{\pi}(s_0) \right) + P(s_1 | s_0, a) \left(r(s_0, a, s_1) + \gamma V_{\pi}(s_1) \right) \right]$$

$$V_{\pi}(s_0) = \sum_{a \in A} \pi(a|s_0) \left(\sum_{s' \in S} P(s'|s_0, a) \left[r(s_0, a, s') + \gamma V_{\pi}(s') \right] \right)$$

Bellman's equation.
for state value function.

* s_0, a_0



$$\underline{q_{\pi}(s_0, a_0)} = \sum_{s' \in S} P(s'|s_0, a_0) \left[r(s_0, a_0, s') + \gamma \sum_{a' \in A} \pi(a'|s') \underline{q_{\pi}(s', a')} \right]$$

$$\left. \begin{aligned} q_*(s_0, a_0) &= 0.7 \\ q_*(s_0, a_1) &= \underline{0.8} \\ q_*(s_0, a_2) &= 0.2 \end{aligned} \right\} \begin{aligned} \pi(a_0|s_0) &= 0 \\ \pi(a_1|s_0) &= 1 \\ \pi(a_2|s_0) &= 0 \end{aligned}$$

$$V_{\pi^*}(s) \longrightarrow \pi^*(a|s)$$

$$V_{\text{OLD}}(s) = 0$$

$$\textcircled{1} \quad V_{\text{NEW}}(s) = \sum_{a \in A} \pi(a|s) \left[\sum_{s' \in S} P(s'|s, a) \left(r(s, a, s') + \gamma V_{\text{OLD}}(s') \right) \right]$$

$$\textcircled{2} \quad \Delta = \|V_{\text{OLD}}(s) - V_{\text{NEW}}(s)\|$$

$$\textcircled{3} \quad V_{\text{OLD}}(s) = V_{\text{NEW}}(s)$$

repeat ① and ②
until $\Delta \leq \text{tol}$

$$\text{tol} = \underline{\underline{0.0001}}$$
