

2.5. Dynamic Programming

(1) Stochastic:

Using eqⁿ 4.4, we can say:

$$\begin{aligned} v_{\lambda}(s) &= E_{\lambda}[G_t | S_t = s] \\ &= \sum_a \lambda(a|s) \sum_{s', r} p(s', r | s, a) [\gamma + \gamma v_{\lambda}(s')] \end{aligned}$$

Using eqⁿ 4.6, we can say

$$\begin{aligned} q_{\lambda}(s, a) &= E[R_{t+1} + \gamma v_{\lambda}(S_{t+1}) | S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r | s, a) [\gamma + \gamma v_{\lambda}(s')] \end{aligned}$$

Using the above two eq^s, we get:-

$$v_{\lambda}(s) = \sum_a \lambda(a|s) q_{\lambda}(s, a).$$

(1) Deterministic :-

$$v^*(s) = Q_{\pi}(s, a) \quad \forall a \in A:$$

$$\pi(a|s) = 1$$

(2) The Bellman optimality eqⁿ for the action-value function:

$$Q_{*}(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} Q_{*}(s', a')]$$

The value-iteration for the state-value function was obtained by turning its Bellman eqⁿ to an update rule. [Barto Sutton]

Using the same for the action-value ~~function~~ iteration:-

$$q_{V_{K+1}}(s) = \sum_{s'} p(s', r | s, a) \left[r + \gamma \max_{a'} q_{V_K}(s', a') \right].$$

(3) We know:

$$v_{\pi}(s) = \sum_{a \in A} \pi(a|s) \cdot q_{\pi}(s, a) \quad \text{--- (1)}$$

$$q_{\pi}(s) = \sum_{s'} p(s', r | s, a) \left[r + \gamma v_{\pi}(s') \right] \quad \text{--- (2)}$$

Replacing $v_{\pi}(s)$ from (1) to (2).

$$q_{\pi}(s) = \sum_{s'} p(s', r | s, a) \left[r + \gamma \sum_{a' \in A} \pi(a'|s) \cdot q_{\pi}(s', a') \right]$$

(4) Using eqⁿ 4.9, we can say

~~is~~

(4) Given:

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$$

We also know that:-

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) \cdot q_{\pi}(s, a) \quad \text{--- (1)}$$

Using (1) we can rewrite the policy improvement step:-

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s', r} p(s', r | s, a) \left[r + \gamma \sum_{a' \in A} \pi(a'|s) \cdot q_{\pi}(s, a') \right]$$

(5) The policy evaluation step on page 75 assumes that $\pi(s)$ has a probability distribution over $a \in A$. That's why a summation over $\forall a \in A$ is taken on page 75.

The evaluation step on page 80 assumes a deterministic or a pure greedy policy with respect to the actions. Hence, there is only one action that should be chosen for a state.