

Reinforcement Learning.

Markov's Decision process:

$$P(S_{t+1} | S_t) = P(S_{t+1} | S_1, \dots, S_t) \text{ - markov}$$

$$P_{ss'} = P(S_{t+1} = s' | S_t = s) \text{ - transition prob.}$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots$$
$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}.$$

Returns.

$$R_s = E[R_{t+1} | S_t].$$

S - states

P - transition prob

R - reward

γ - discount factor.

$$\pi(a|s) = P(A_t = a | S_t = s)$$

policy fn.

$$\text{state value fn} = V_{\pi}(s) = E_{\pi}[G_t | S_t = s].$$

(for a given policy).

$$V(s) = E[R_{t+1} + \gamma V(S_{t+1}) | S_t = s].$$

Bellman.

State-action Value fn.

$$q_{\pi}(s, a) = E_{\pi} [G_t | S_t = s, a_t = a]$$
$$= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, a_t = a \right]$$

$$V(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} V(s')$$

$$\left[\begin{array}{l} V = R_s + \gamma P V \\ (1 - \gamma P)V = R_s \\ V = (1 - \gamma P)^{-1} R_s. \end{array} \right. \leftarrow \text{Bellman linear eq.}$$

Bellman expectation eq:

$$V(s) = E [R_{t+1} + \gamma V(S_{t+1}) | S_t = s].$$