

MARKOV DECISION PROCESSES

dynamic setting

v and q -values

Tentative taxonom

RL algorithms

MDP:

states : S set of states

S finite or S infinite

actions : A set of actions

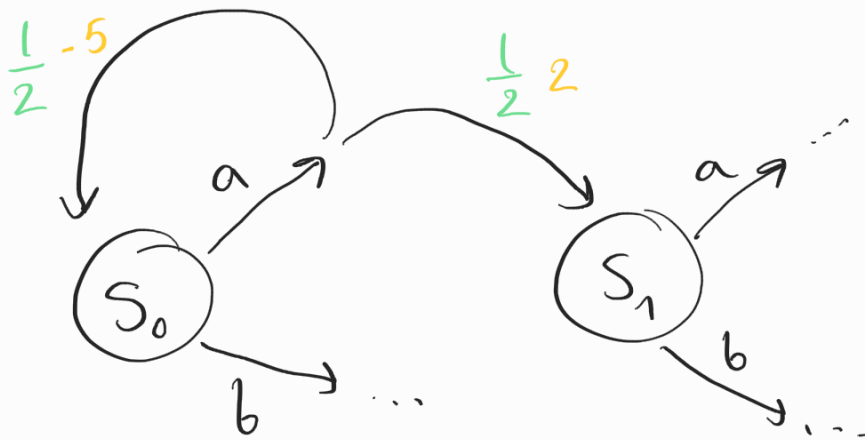
set of reals

A finite or A infinite

probabilistic transition function: $\Delta: S \times A \rightarrow \text{Dist}(S \times \mathbb{R})$

rewards

$\Delta(s, a)(s', r)$: probability that from state s playing action a , we go to state s' and get reward r



MDP : S , A , Δ

states , actions , transition function

Goal: Construct a strategy / policy

$$\sigma: S \rightarrow A$$

such that σ maximises

$$\lambda \in (0, 1)$$

$$\sigma \left[\sum_{t=0}^{\infty} \lambda^t r_t \right]$$

$r_0 + \lambda r_1 + \lambda^2 r_2 + \lambda^3 r_3 + \dots$

$s_0 \in S$ initial state

$$\sigma: \begin{array}{l} \sigma(s_0) = a_0 \\ \sigma(s_1) = a_1 \end{array} \quad \begin{array}{l} \Delta(s_0, a_0)(s_1, r_0) \\ \Delta(s_1, a_1)(s_2, r_1) \end{array}$$

$s_0, a_0, r_0, s_1, a_1, r_1, s_2, \dots$

play / trajectory / path
infinite or finite

Two cases:

- either eventually we reach a sink

$$\sum_{t=0}^{\infty} r_t \text{ is actually finite}$$

→ FINITE HORIZON

- or the trajectory may be infinite

$$\sum_{t=0}^{\infty}$$

$\sum_{t=0}^{\infty} \gamma^t r_t$ may not be defined

→ DISCOUNTED

$\gamma \in (0, 1)$: fixed constant

$$\sum_{t=0}^{\infty} \gamma^t r_t = r_0 + \underbrace{\gamma r_1}_{\gamma^1 r_1} + \underbrace{\gamma^2 r_2}_{\gamma^2 r_2} + \gamma^3 r_3 + \dots$$

$\gamma^t \xrightarrow{t \rightarrow \infty} 0$

Good news:

- (1) well defined mathematically
- (2) reasonable in practice