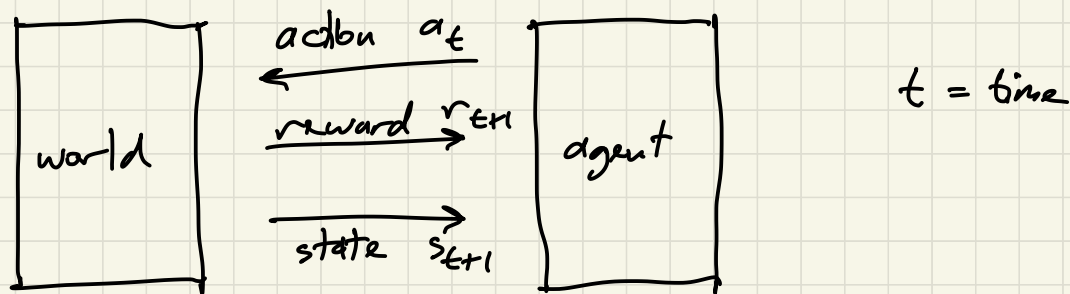


MDPs - Markov Decision Processes

- stochastic, discrete state, discrete action, state feedback (fully observable)

environment



3 system spaces: state $s_t \in S$, action $a_t \in A$, reward $r_t \in \mathbb{R}$

3 functions: model $p(s_{t+1} | s_t, a_t) = \text{prob. of } s_{t+1} \text{ given } s_t, a_t$

state transition
function

reward $r_{t+1} = R(s_t, a_t) \leftarrow \text{deterministic}$

we want to find \rightarrow policy $\pi(a_t | s_t) = \text{prob. of } a_t \text{ given } s_t$

Notation

lowercase $s, a, r, s_t, a_t \leftarrow$ values, samples, realizations

uppercase $S_t, A_t, R_t \leftarrow$ random variables

$$p(s_{t+1} | s_t, a_t) = \text{Prob}(S_{t+1} = s_{t+1} | S_t = s_t, A_t = a_t)$$

$$p : S \times S \times A \rightarrow \mathbb{R}$$

$s_{t+1} \sim p(\cdot | s_t, a_t) \leftarrow s_{t+1}$ has (marginal) distribution given by p
 s_{t+1} is sampled from p

$S \quad S \quad A$
 $\downarrow \downarrow \downarrow$
def $p(s_2, s_1, a_1) =$

return r

\mathbb{R}

Example $f(s) =$ expected (average) 1-step reward from state s

$$= E[R(s_t, A_t) | s_t = s] = E[R(s, a) | s]$$

\uparrow R.V.

$$= E_{a \sim \pi(\cdot | s)} [R(s, a)]$$

$$= \sum_{a \in A} \pi(a | s) R(s, a)$$