

# Markov Decision Process (with reward)

$$\mathcal{M} = (S, \mathcal{M}_0, A, P, R, \gamma)$$

$S$ : state space, set of all possible states

$\mathcal{M}_0$ : initial state distribution;  $\mathcal{M}_0 \in \Delta(S)$

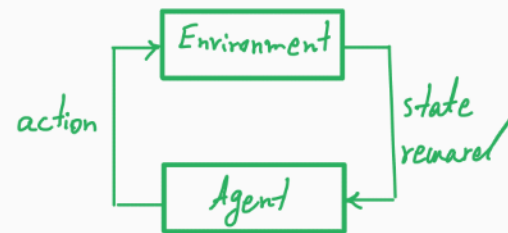
$A$ : action space, set of all possible actions

$P$ : transition function;  $P: S \times A \longrightarrow \Delta(S)$

$R$ : reward function;  $R: S \times A \longrightarrow \mathbb{R}$

$\gamma$ : discount factor;  $\gamma \in [0, 1]$

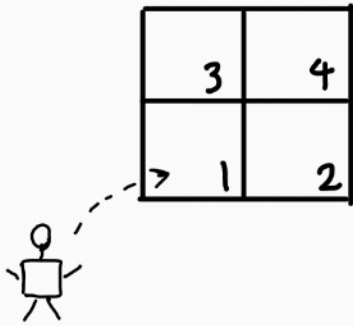
Agent's interaction with the environment:



- The interaction starts at time step  $t=0$  and state  $s_0 \sim \mathcal{M}_0$ .
- At time step  $t \in \mathbb{N}_0$  and state  $s_t \in S$ , the agent takes action  $A_t \in A$ , observes the next state  $s_{t+1} \sim P(\cdot | s_t, A_t)$ , and receives the immediate reward  $R_t = R(s_t, A_t) \in \mathbb{R}$ .
- The history of interaction at time  $t$  is called a trajectory  $\tau_t \in \mathcal{T}$ , where

$$\tau_t = (\underbrace{s_0, a_0, r_0}_{t=0}, \underbrace{s_1, a_1, r_1}_{t=1}, \dots, \underbrace{s_t, a_t, r_t}_{t=t})$$

Example:



$$S = \{1, 2, 3, 4\}$$

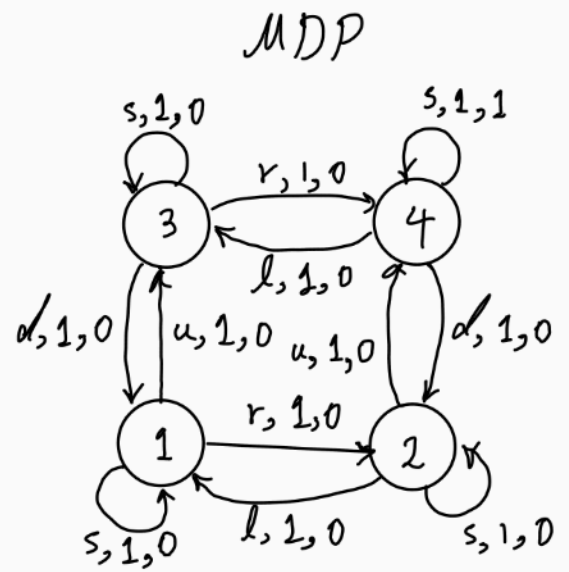
$$A = \{r, l, u, d, s\}$$

$$M_0 = [1, 0, 0, 0]$$

sample trajectories

( $s_0=1, a_0=u, r_0=0, s_1=3, a_1=s, r_1=0, s_2=3, \dots$ ) edge label: (action, transition probability, reward)

( $s_0=1, a_0=r, r_0=0, s_1=2, a_1=u, r_1=0, s_2=4, a_2=s, r_2=1, \dots$ )



**Policy:** A possibly randomized mapping from a trajectory to actions;

$$\pi: \mathcal{T} \rightarrow \Delta(A)$$

set of all possible trajectories

$(s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{H-1}, a_{H-1}, r_{H-1})$

$n$  possibilities  $\downarrow$  deterministic  $\downarrow$   $n \times m$  possibilities

- stationary policy: A policy where the action only depends on the current state;

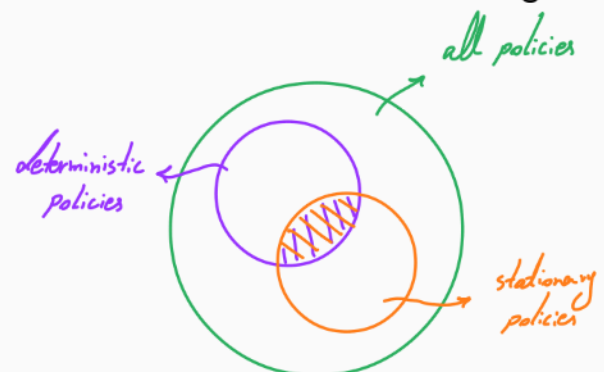
$$\pi: S \rightarrow \Delta(A)$$

- Deterministic policy: A policy where the action is chosen deterministically;

$$\pi: \mathcal{T} \rightarrow A$$

- Deterministic, stationary policy:

$$\pi: S \rightarrow A$$



**Goal:** The learning agent aims to find a policy  $\pi$  that maximizes the expected (discounted) cumulative reward

$$\pi^* \in \max_{\pi \in \Pi} \mathbb{E}_{\substack{S_0 \sim \mu_0 \\ A_t \sim \pi(\tau_t) \\ S_{t+1} \sim P(\cdot | S_t, A_t)}} \left[ \overbrace{\sum_{t=0}^T \gamma^t R(S_t, A_t)}^{\text{return}} \right]$$

for stochastic reward  $\rightarrow R(S_t, A_t) \sim \mathcal{R}(\cdot, \cdot)$

$T$   $\rightarrow$   $T = \infty$  infinite horizon  
 $T$   $\rightarrow$   $T < \infty$  finite horizon

$\gamma$   $\rightarrow$   $\gamma = 1$  undiscounted  
 $\gamma$   $\rightarrow$   $0 < \gamma < 1$  discounted  
 $\gamma$   $\rightarrow$   $\gamma = 0$  greedy / myopic

\* We will start with the infinite-horizon discounted setting.

## Value functions

- (state) value function of a policy  $\pi$ ,  $V^\pi: S \rightarrow \mathbb{R}$ , is defined as

$$V^\pi(s) = \mathbb{E}_{\substack{A_t \sim \pi(\tau_t) \\ S_{t+1} \sim P(\cdot | S_t, A_t)}} \left[ \sum_{t=0}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s \right] \quad \forall s \in S$$

- (state-) action value function of a policy  $\pi$ ,  $Q^\pi: S \times A \rightarrow \mathbb{R}$  is defined as

$$Q^\pi(s, a) = \mathbb{E}_{\substack{A_t \sim \pi(\tau_t) \\ S_{t+1} \sim P(\cdot | S_t, A_t)}} \left[ \sum_{t=0}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s, A_0 = a \right] \quad \forall s \in S$$

$t \geq 1$   $\leftarrow$

## Bellman (consistency) Equations:

Let  $\pi$  be a stationary policy. The value functions under  $\pi$ , i.e.,  $V^\pi$  and  $Q^\pi$ , satisfy the following equations:

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(s)} \left[ \overbrace{R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [V^\pi(s')]}^{Q^\pi(s,a)} \right]$$

$$Q^\pi(s,a) = R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [V^\pi(s')]$$

$$Q^\pi(s,a) = \mathbb{E}_{\substack{A_t \sim \pi(s) \\ t \geq 1 \leftarrow s_{t+1} \sim P(\cdot|s_t, A_t)}} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, A_t) \mid s_0 = s, A_0 = a \right] \quad \forall s \in S$$

$$= \mathbb{E}_{\substack{A_t \\ s_{t+1}}} \left[ \underbrace{R(s_0, A_0)}_{R(s,a)} + \sum_{t=1}^{\infty} \gamma^t R(s_t, A_t) \mid s_0 = s, A_0 = a \right]$$

$$= R(s,a) + \gamma \mathbb{E}_{\substack{A_t \\ s_{t+1}}} \left[ \underbrace{\sum_{t=1}^{\infty} \gamma^{t-1} R(s_t, A_t) \mid s_0 = s, A_0 = a}_{\left[ \sum_{t'=0}^{\infty} \gamma^{t'} R(s_{t'+1}, A_{t'+1}) \mid s_0 = s, A_0 = a \right]} \right]$$

$$\mathbb{E}_{s' \sim P(\cdot|s,a)} \left[ \underbrace{\sum_{t'=0}^{\infty} \gamma^{t'} R(s_{t'+1}, A_{t'+1}) \mid s_{t'=0} = s'}_{V^\pi(s')} \right]$$