# Markov Decision Processes

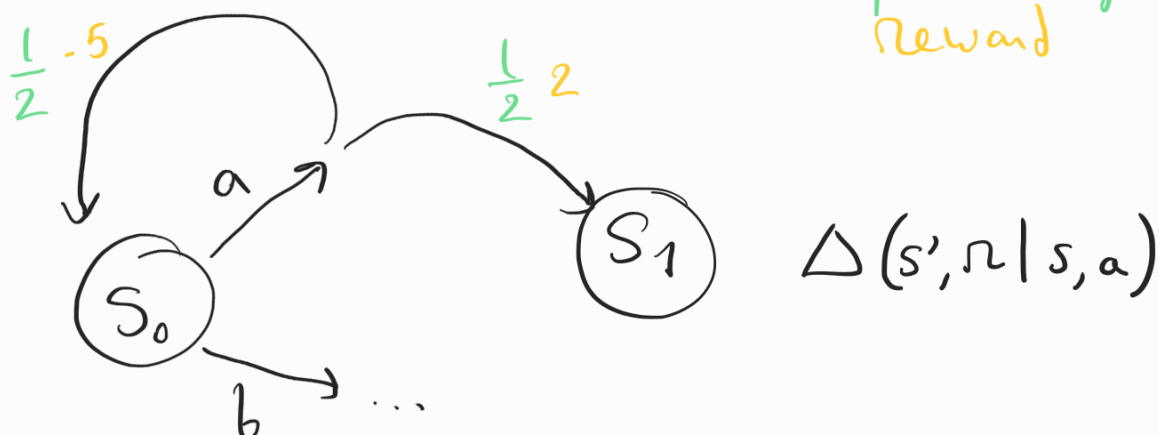**states** : $S$ set of states

**actions** : $A$ set of actions

**transition function** : $\Delta : S \times A \longrightarrow Dist(S \times \mathbb{R})$

rewards

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

probability
reward



$\frac{1}{2}$ .5  $\frac{1}{2}$ 2

$a$   $S_1$   $\Delta(s', r \mid s, a)$

$S_0$

$b$   ...

**Strategy ≡ policy :**

$$\pi : S \longrightarrow A \qquad \text{deterministic}$$

on $\quad \pi : S \longrightarrow \text{Dist}(A) \qquad$ stochastic

distributions

**play ≡ trajectory ≡ path :**

$$\rho = S(0), A(0), R(0), S(1), A(1), R(1), \ldots$$

**return of a trajectory**

$$G = \sum_{t \geq 0} \gamma^t R(t) = R(0) + \gamma R(1) + \gamma^2 R(2) + \ldots$$

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

Two cases:

- either eventually we reach a sink

$$G = \sum_{t=0}^{\infty} R(t) \quad \text{is actually finite}$$

$\longrightarrow$ FINITE HORIZON

- or the trajectory may be infinite

$$G = \sum_{t=0}^{\infty} \gamma^t R(t)$$

$\longrightarrow$ DISCOUNTED

$$\gamma \in (0,1) \quad : \text{fixed constant}$$

$$G = \sum_{t=0}^{\infty} \gamma^t R(t) = R(0) + \gamma R(1) + \gamma^2 R(2) + \gamma^3 R(3) \dots$$

$$\gamma^t \xrightarrow[t \to \infty]{} 0$$

**Goal:** Construct a strategy $\Pi$

maximising $\mathbb{E}\left[ G \mid S_{(0)} = S_0 \wedge \Pi \right]$

$S_0 \in S$      initial state