

MDP

- A mathematical formulation of the environment and what we want the agent to learn.

- $t \in \{0, 1, \dots\}$ time step

- $S_t, A_t, R_t \leftarrow$ Reward given to the agent at time t .
 \uparrow state at time t \uparrow Action chosen at time t

MDP $M = (S, A, P, R, \gamma)$

1. $S =$ Set of all possible states of the env. [finite]
2. $A =$ Set of all possible actions. [Finite]
3. $P =$ "Transition function" describes how state of the env. transition.

Finite assumption for MDP (as opposed to continuous)

$$P: S \times A \times S \Rightarrow [0, 1]$$

$$p(s, a, s') \triangleq P_x(S_{t+1} = s' | S_t = s, A_t = a)$$

\uparrow
[is defined to be]

$$\forall s \in S, a \in A, s' \in S$$

4. $R:$ "reward function" describes how rewards are generated. $R: S \times A \times S \rightarrow \mathbb{R}$.
 $R(s, a, s') = \mathbb{E}[R_t | S_t = s, A_t = a, S_{t+1} = s']$

1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	25

$$p(8, \text{right}, 9) = 0.8$$

$$p(8, \text{right}, 15) = 0$$

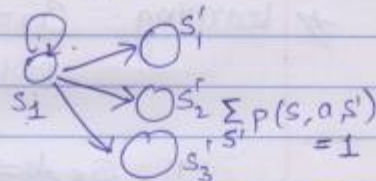
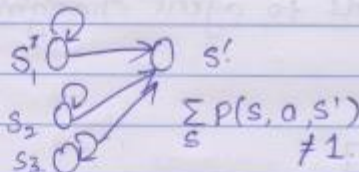
Deterministic $\Rightarrow p(s, a, s') \in \{0, 1\}$

R: "Reward" $R(20, \text{right}, 1) \rightarrow \text{undefined}$

$$R(20, \text{right}, 20) \Rightarrow 0$$

$$R(20, \text{right}, 21) = -10 \rightarrow R_t$$

What is $\sum_{s'} p(s, a, s')$ and $\sum_s p(s, a, s')$?



A note on R_t :

$$R_t(20, \text{right}, 21) = -10 \quad \text{[We're already conditional on } s_{t+1} \text{ (see formula) for } R_t]$$

$$R(20, \text{right}) = 0.8 \times -10 + 0.05 \times 0 + 0.5 \times 0 + 0.2 \times 0 + \dots$$

used answering a question. Hereafter we use $R(s, a, s')$, not $R(s, a)$

Why take an expectation over R ?

5. d_0 : Initial state distribution

$$d_0: S \rightarrow [0, 1]$$

$$d_0(s) \triangleq P_0(s_0 = s)$$

6. γ = Reward discount parameter

$$\gamma \in [0, 1]$$

Agent formulation (for an MDP)

★ Policy: The mechanism within the agent that determines which action to take in a state.

★ Learning: Corresponds to agent changing its policy.

$$\pi: S \times A \Rightarrow [0, 1]$$

$$\pi(s, a) \triangleq \Pr(A_t = a | S_t = s)$$

★ Deterministic policies: A policy is deterministic if it always chooses same action in a given state.

$$\pi(s, a) \in \{0, 1\} \dots \pi \text{ is deterministic}$$

★ Stochastic policies:

π	actions			
	up	down	left	right
1	0.01	0.01	0.08	0.9
2	0	0	0.5	0.5
3			1	
4			1	
state 5				1

All matrices with positive entries and rows summing to 1 are policies

23
 (Representation of π)

$M = (S, A, P, R, \gamma)$

$s_0 \sim d_0$ (initial state sampled)

$a_0 \sim \pi(s_0, \cdot)$ (a_0 sampled)

$s_1 \sim p(s_0, a_0, \cdot)$ notation to denote "only final state"

R_0 (random variable) \rightarrow Computed with $\mathbb{E}[R_0] = R(s_0, a_0, s_1)$

\hookrightarrow Deterministic reward

$R_0 = R(s_0, a_0, s_1)$

★ Agent's goal:

Find a policy that maximizes the expected amount of reward that it will get.

Objective function: $J: \Pi \xrightarrow{\text{set of all policies}} \mathbb{R}$

$$J(\pi) \triangleq \mathbb{E} \left[\sum_{t=0}^{\infty} R_t \mid \pi \right]$$

(This is a common objective function, but there are others.)

↑
But a policy isn't an event, so what does it mean to condition on π ?

- We assume a fixed policy (π) that ~~is~~ used across all time steps.

$$\begin{aligned} \hookrightarrow A_0 &\sim \pi(s_0, \cdot) \\ A_1 &\sim \pi(s_1, \cdot) \\ A_2 &\sim \pi(s_2, \cdot) \end{aligned}$$

- This objective treats rewards as additive quantities, but other objective functions ~~could~~ ^{may} exist.
not. $\&$

Optimal policy

$$\pi^* \in \operatorname{argmax}_{\pi \in \Pi} J(\pi) \quad \left[\operatorname{argmax} \text{ returns the set of optimal policies (there can be multiple)} \right]$$

If $|S|$ and $|A|$ are finite, and R_t is bounded, then an optimal policy exists.