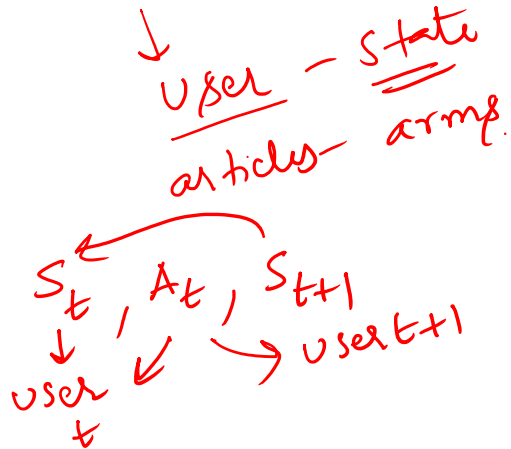


Markov Decision Processes (MDP)

Prof. Subrahmanya Swamy



Multi-Arm

- one state ✓

Contextual Bandits ✓

- News article

RL Framework

2

$S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_T$
 $\rightarrow "R_1 + R_2 + \dots + R_T"$

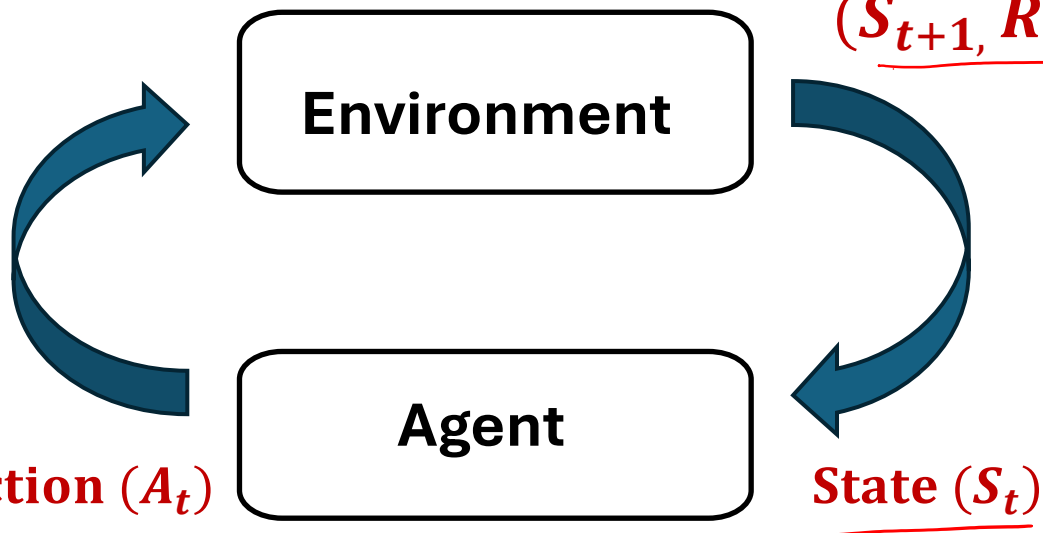
Next State, Reward
 (S_{t+1}, R_{t+1})

1. Agent observes the state ✓
2. Takes an action ✓
3. Environment puts the agent in a new state &
4. Also gives a reward based on taken action

Goal:

Learn policy to maximize the cumulative reward $\sum_t R_t$

$$\pi : S \rightarrow A$$



How do we mathematically model the State transitions and Rewards?

MDP

Independent Random Variables

- A sequence of coin tosses X_1, X_2, X_3, \dots
 \downarrow
 $1/0$
- Head: 1, Tail: 0, Bias of coin: p_h \rightarrow Prob (Head)
 $\underline{p_h} \quad P(X_1=1) = p_h$
- Knowledge of X_1 does not help in predicting X_2
 \downarrow
 x_n
- $\mathbb{P}(X_2 = 1 \mid X_1 = 0) = p_h$ ✓
 \downarrow
- $\mathbb{P}(X_2 = 1 \mid X_1 = 1) = p_h$ ✓
 \uparrow

Markov Chain ✓

- A sequence of coin tosses X_1, X_2, X_3, \dots

↓

Independent

1/-1

$X_1 = 1$ $X_2 = 1$ $X_3 = -1$

$Y_0 = 0$ $Y_2 = 1 + 1 = 2$ $Y_3 = X_3 + Y_2$

$Y_1, Y_2, \dots, Y_t, \dots$

↑ ↑ ↑

-1 2 1

- If coin lands in

- Head: Win 1 rupee ✓
 - Tail: Lose 1 rupee ✓
- 1 Lose

- Define Y_t = total money accumulated till time t

- Y_1, Y_2, Y_3, \dots are dependant RVs

- $\mathbb{P}(Y_5 = 1 | Y_4 = 3) = 0$
 - $\mathbb{P}(Y_5 = 1 | Y_4 = 0) = \frac{1}{2}$
- ↑

→

$3 + 1 = 4$

$3 - 1 = 2$

Markov Chain

$$Y_1, \dots, Y_t$$

$$S_{\text{states}} = \{0, 1, 2, \dots\}$$

- Y_1, Y_2, Y_3, \dots satisfy Markov property!

- **Markov Property:** Given the present, the future is independent of the past!

$$\mathbb{P}(Y_5 = 1 | Y_4 = 2, Y_3 = 3) = \frac{1}{2} \quad \checkmark \quad P(\text{Tail})$$

$$\mathbb{P}(Y_5 = 1 | Y_4 = 2, Y_3 = 1) = \frac{1}{2} \quad \checkmark$$

$$t=4$$

$$-1, -1 \rightarrow 3-2=1$$

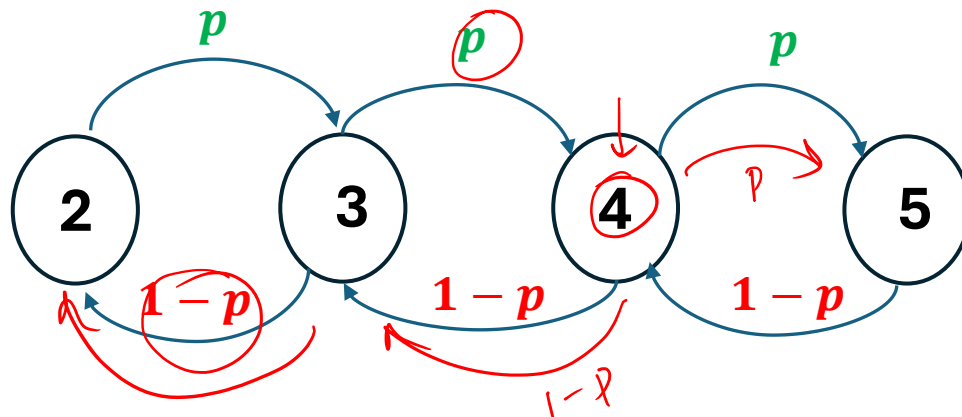
$$+1, -1 \rightarrow 1-2=-1$$

$$Y_3=3$$

$$Y_3=1$$

$$Y_4=2$$

$$Y_5=1$$



Markov Chain Specification (S , $P_{ss'}$)

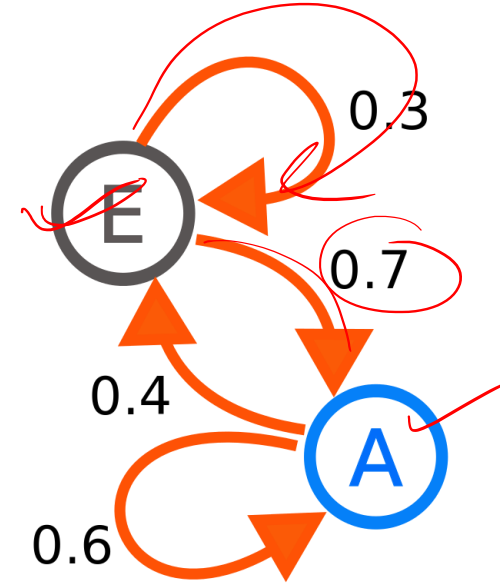
• $S \rightarrow$ State space $\{E, A\}$ MDP

• $P_{ss'}$ \rightarrow Transition probability

$$\mathbb{P}(\underline{S_{t+1}} = s' \mid \underline{S_t} = s)$$

	<u>E</u>	<u>A</u>
<u>E</u>	0.3	0.7
<u>A</u>	0.4	0.6

$\underline{s}, \underline{s'} \in S$
 $\underline{s_t}, \underline{s_{t+1}}$
 $\underline{A_t}$



Markov Decision Process (MDP)

- Introduce action to convert Markov Chain into MDP

S_t

- Actions: How much money to bet (A_t) in the game when I have Y_t money?

- If $Y_t = 3$, then possible actions are $\{1, 2, 3\}$.

$S = \{1, 2, \dots\}$

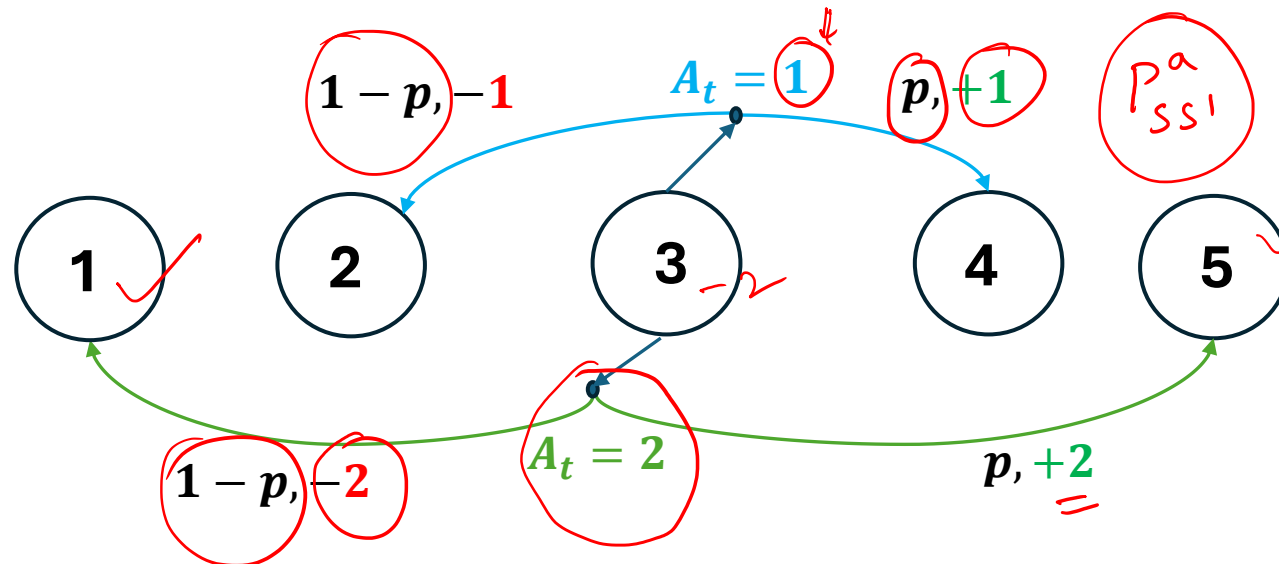
$A = \{$

R

$P^1_{3,4} = p$

$P^2_{3,1} = 1-p$

$P^2_{3,6} = 0$

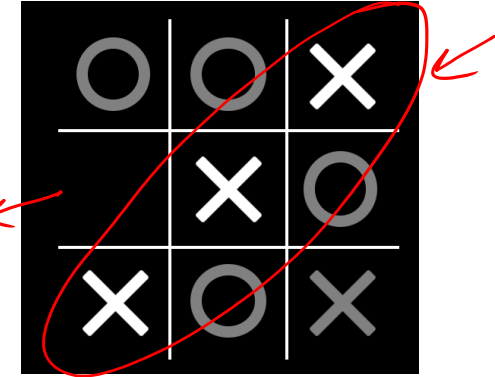


Episodic and Continuing MDPs

• Episodic ✓ Task.

→ $s_0, a_0, r_1, s_1, \dots, s_T$

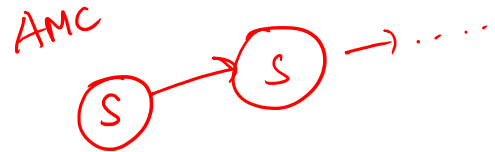
Terminal state in Tic-Tac-Toe



- There **exists** a special state called the terminal state
- The episode ends at the terminal state
- Eg: Board games

$s_0, s_1, s_2, s_3, \dots$

• "Continuous" Task.



- No terminal state exists
- The task continues forever
- Eg: Portfolio management ✓
 - Every day, decide which shares to buy/sell

Discount Factor in MDP

- Episodic task:

- Total Reward (Return) : $G_t = R_{t+1} + R_{t+2} + \dots + R_T$ ✓
 \downarrow
 $\leq M = 1$
- Bounded Returns if each $R_i \leq M$
 \downarrow
 $\leq M$

- Continuing task: ✓

- $G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots$
- $G_t = \sum_{i=t+1}^{\infty} R_i$ could become unbounded even if each $R_i \leq M$
 \downarrow
 $\leq M$

- Solution: Discount factor $\gamma \in (0, 1)$

- $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$ $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$
- $G_t = \sum_{i=0}^{\infty} \gamma^i R_{t+1+i} \leq \frac{M}{1-\gamma}$ (Bounded)
 \downarrow
 $\leq \frac{M}{1-\gamma}$
 \downarrow
 $\gamma = 0.9$
 $\gamma = 0.99$
 $\gamma = 0.001$
- High $\gamma \sim 1 \Rightarrow$ Long-term planning ✓
- Low $\gamma \sim 0 \Rightarrow$ Short-term planning

MDP Specification ($\underline{S}, \underline{A}, \underline{R}_S^a, \underline{P}_{SS'}^a, \underline{\gamma}$)

• $\underline{S} \rightarrow$ State space (incl. terminal states if any)

• $\underline{A} \rightarrow$ Action space

• $\underline{R}_S^a \rightarrow$ Expected Rewards

• $\mathbb{E}[R_{t+1} \mid S_t = s, A_t = a]$

• $\underline{P}_{SS'}^a \rightarrow$ Transition probabilities

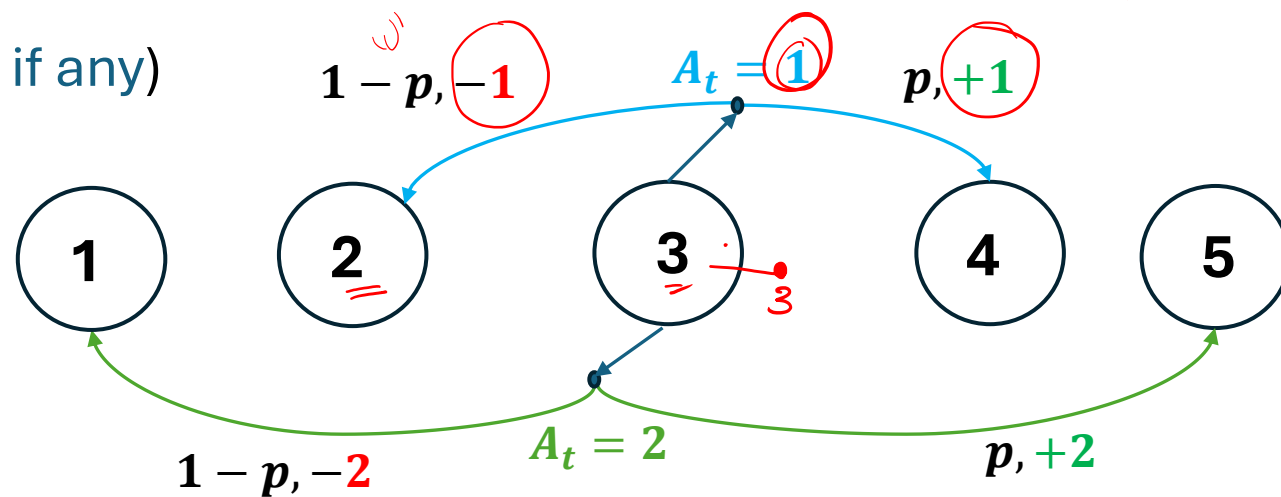
• $\mathbb{P}(S_{t+1} = s' \mid S_t = s, A_t = a)$

• $\underline{\gamma} \in (0,1) \rightarrow$ Discount factor

$$E[x] = \sum_x p_x \cdot x$$

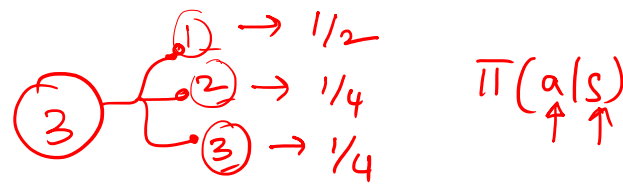
$$R_s^a$$

$$\begin{aligned} R_3^1 &= +1(P) - 1(1-P) \\ &= P - 1 + P \\ &= 2P - 1 \end{aligned} \quad 10$$



Optimal Policy

11

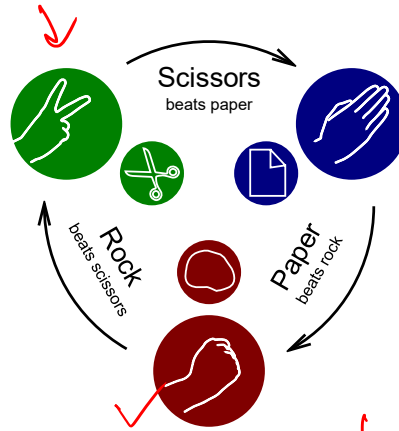


- Policy:

- **Deterministic:** $\pi(s): \mathcal{S} \rightarrow \mathcal{A}$ Which action to take in state s
- **Stochastic:** $\pi(a | s)$ In state s , with what probability to take action a

- Why stochastic policies?

- Partially observed states
- Exploration

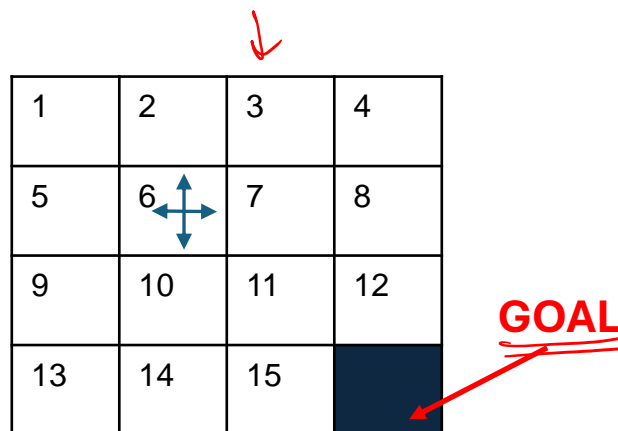
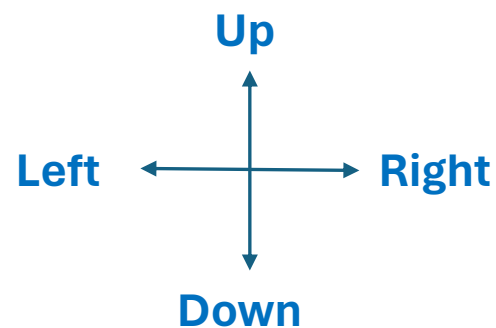


- Optimal Policy:

- π that maximizes the expected return $\mathbb{E}_{\pi}[G_t | S_t = s]$ from any state s

How to model your problem as an MDP?

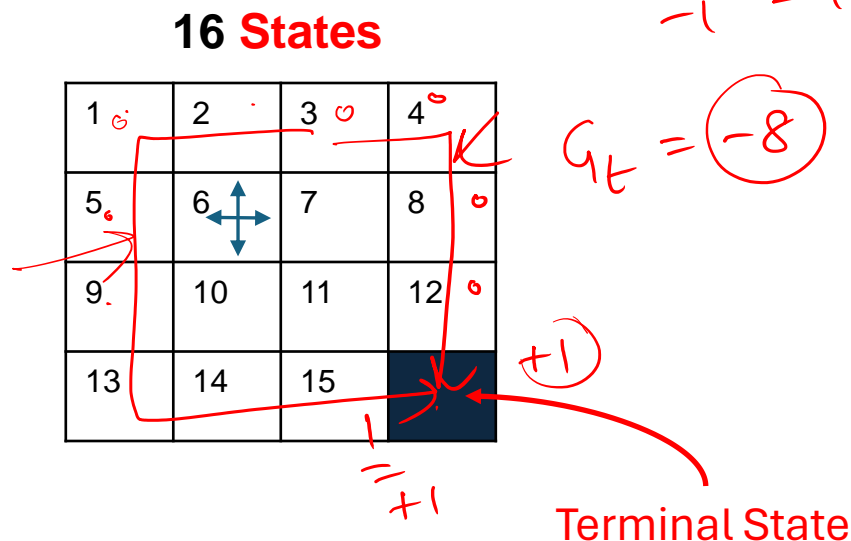
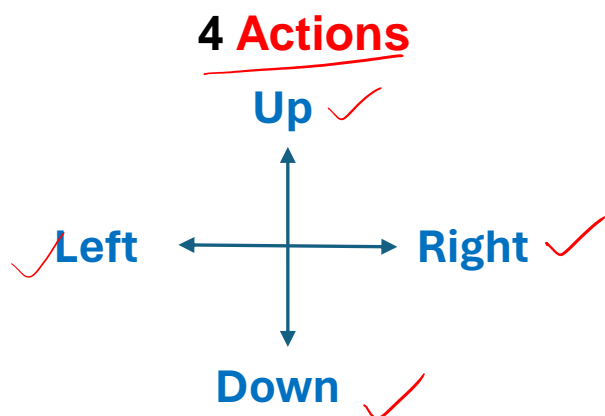
Maze Solving Problem: To reach the goal in the shortest path!



• How to formulate this maze-solving problem as an MDP?

- States ? ✓
- Actions ? ✓
- Rewards ? ✓ R_s^a
- Transition Probabilities ? ✓ $P_{ss'}^a$
- Discount factor ? ✓ γ

How to model your problem as an MDP?



Rewards

$R_t = -1$ on all transitions

Discount Factor

$\gamma = 1$

Terminal State

Deterministic State transitions : $\mathbb{P}(S_{t+1} = 2 \mid S_t = 6, A_t = Up) = 1$

Verify that optimal policy = shortest path

$$E_{\pi}[G_t \mid s_t = s]$$

Exercise ✓

- Alternate MDP formulation for the Maze problem
- Instead of giving -1 reward per each step, can we give 0 reward for every action except for the final action that leads us to the Goal State?
↓ ↓
↑ ↑
- Does the optimal policy of this alternate MDP learn the shortest path?
↓
?
.
- **Hint:** What "discount factor" will help here?

$$\gamma =$$