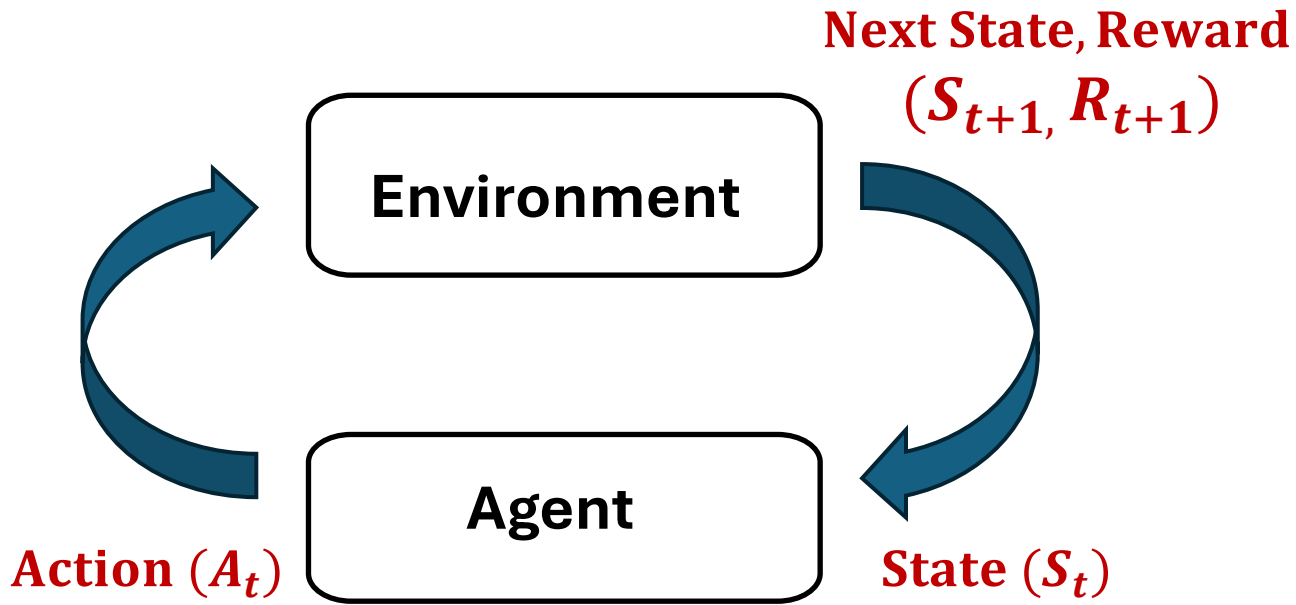


Markov Decision Processes (MDP)

Prof. Subrahmanya Swamy

3

RL Framework



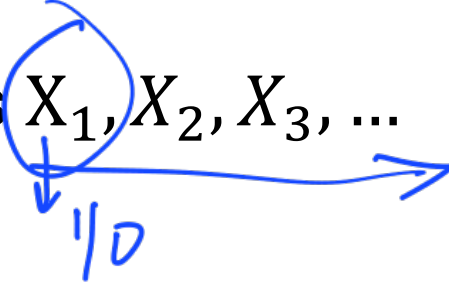

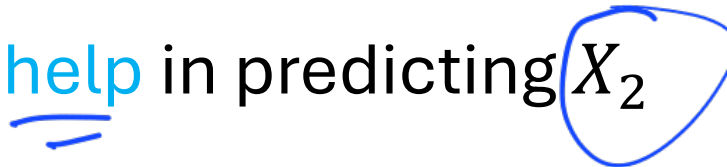
1. Agent observes the state
2. Takes an action
3. Environment puts the agent in a new state &
4. Also gives a reward based on taken action

Goal:

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

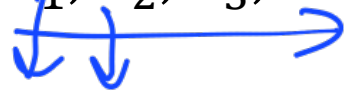
How do we mathematically model the State transitions and Rewards?

Independent Random Variables

- A sequence of coin tosses X_1, X_2, X_3, \dots 
- Head: 1, Tail: 0, Bias of coin: p_h 
- Knowledge of X_1 does not help in predicting X_2 
- $\mathbb{P}(X_2 = 1 | X_1 = 0) = p_h$
- $\mathbb{P}(X_2 = 1 | X_1 = 1) = p_h$

Markov Chain

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



- If coin lands in

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

$$Y_1 \quad Y_2$$

- Define Y_t = total money accumulated till time t

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

$$Y_{t+1} = Y_t + X_t$$

↑
+1 / -1

$$\mathbb{P}(Y_5 = 1 | Y_4 = 3) = 0$$

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

Markov Chain

$$P(Y_5 | Y_3, Y_4)$$

$$P(Y_5 | \overset{\downarrow}{Y_4}, Y_3, Y_2, \overset{5}{Y_1})$$

$$P(Y_5 | Y_4)$$

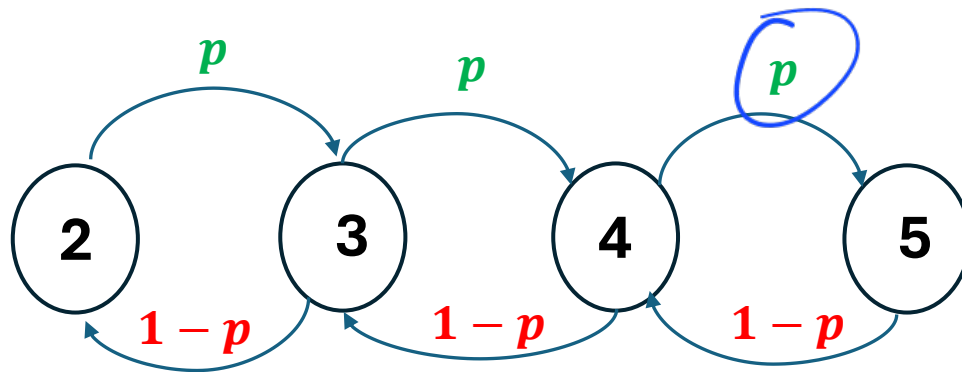
- Y_1, Y_2, Y_3, \dots satisfy Markov property!

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

- $\mathbb{P}(Y_5 = 1 | Y_4 = 2, Y_3 = 3) = \frac{1}{2}$
- $\mathbb{P}(Y_5 = 1 | Y_4 = 2, Y_3 = 1) = \frac{1}{2}$

$$\begin{matrix} 1 & -1 \\ -1 & -1 \end{matrix} \quad \begin{matrix} 1/2 \\ 1/2 \end{matrix}$$

$$P(Y_5 = 0 | Y_4 = 2) = 0$$



$$P = 1/2$$

$$\begin{matrix} \uparrow & \downarrow \\ 3 & 1 \\ \hline P(H) & P(T) \\ p & 1-p \end{matrix}$$

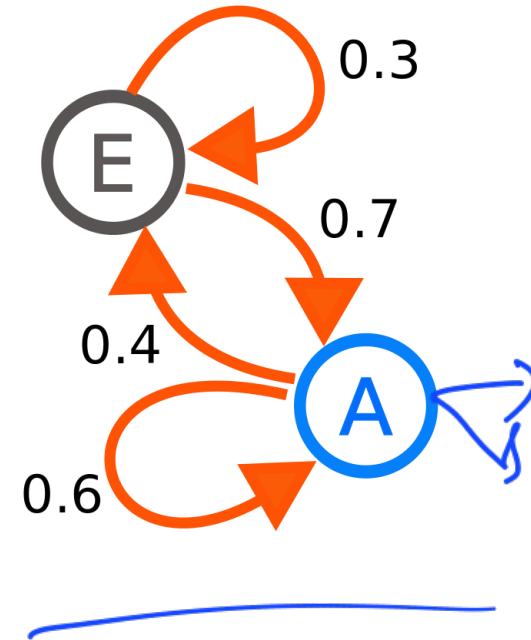
Markov Chain Specification $(S, P_{ss'})$

- $S \rightarrow$ State space $\{E, A\}$
-

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

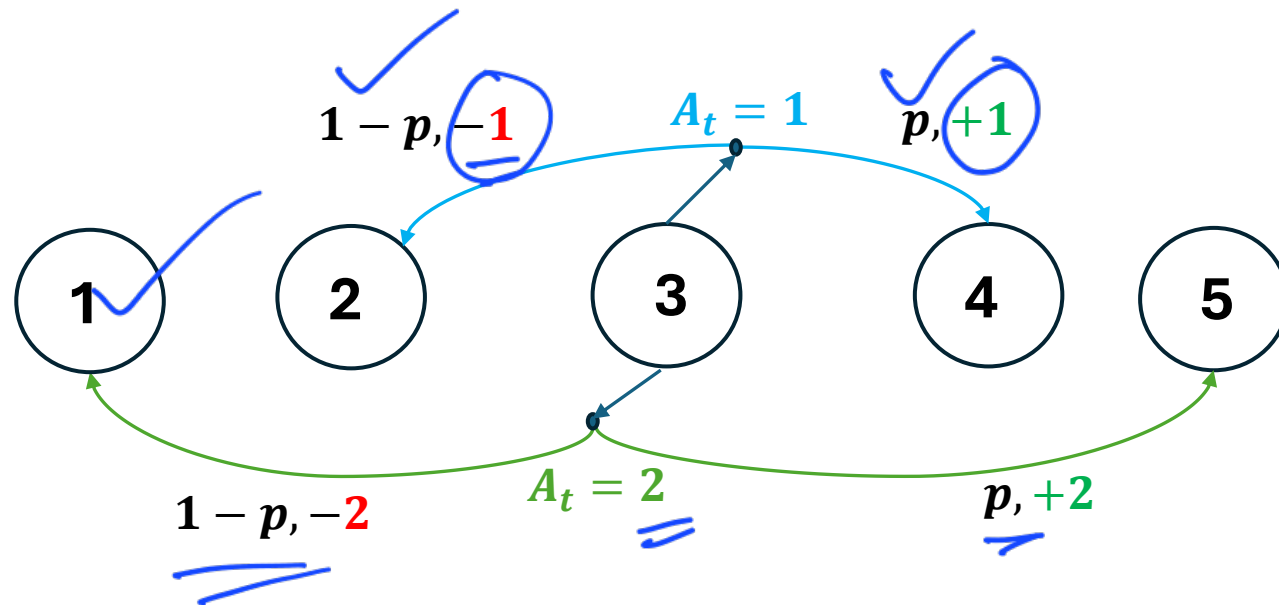
- $\mathbb{P}(S_{t+1} = s' \mid S_t = s)$

	E	A
E	0.3	0.7
A	0.4	0.6



Markov Decision Process (MDP)

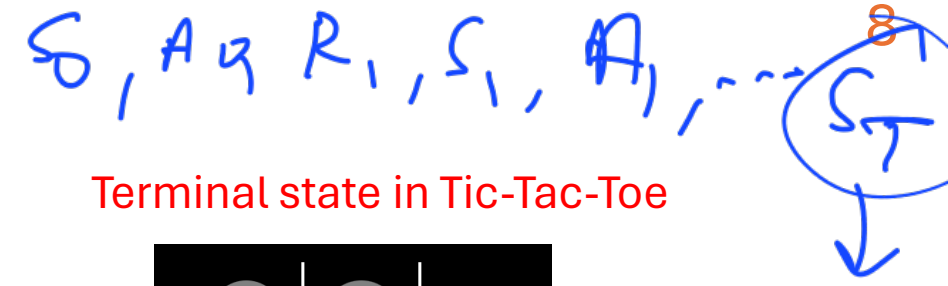
- Introduce action to convert Markov Chain into MDP
- Actions: How much money to bet (A_t) in the game when I have Y_t money?
- If $Y_t = 3$, then possible actions are $\{1, 2, 3\}$.



$$R_s^a = E[R_{t+1} | S_t = s, A_t = a]$$

$$P(S_{t+1} = s' | S_t = s, A_t = a)$$

Episodic and Continuing MDPs



• Episodic

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

• Continuous

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

S_0, A_0, \dots
AMC

Discount Factor in MDP

- Episodic task:

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

- 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 $R_i \leq M$

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

- $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$
- $G_t = \sum_{i=0}^{\infty} \gamma^i R_{t+1+i} \leq \frac{M}{1-\gamma}$ (**Bounded**)

- High $\gamma \sim 1 \Rightarrow$ Long-term planning

- Low $\gamma \sim 0 \Rightarrow$ Short-term planning

$$\leq M + \gamma M + \gamma^2 M + \dots = \frac{M}{1-\gamma}$$

Handwritten example: $\frac{1}{2} \left(\frac{1}{4}\right)^2$ with an arrow pointing to the exponent 2.

MDP Specification $(S, A, R_s^a, P_{ss'}^a, \gamma)$

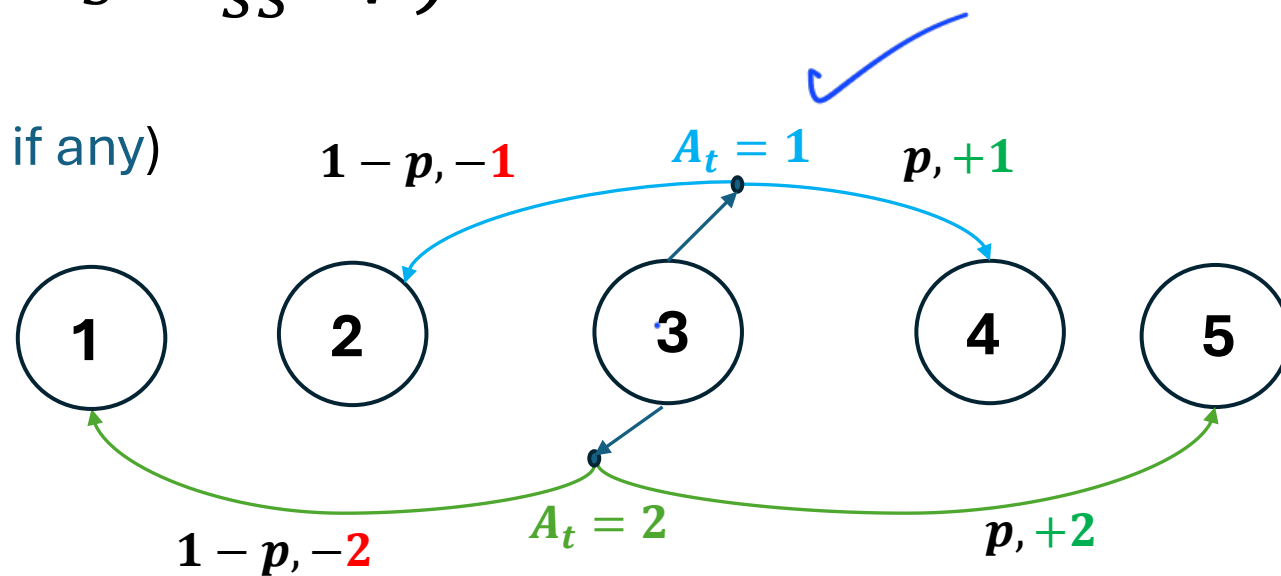
- $S \rightarrow$ State space (incl. terminal states if any)
- $A \rightarrow$ Action space

- $R_s^a \rightarrow$ Expected Rewards
- $\mathbb{E}[R_{t+1} \mid S_t = s, A_t = a]$

- $P_{ss'}^a \rightarrow$ Transition probabilities
- $\mathbb{P}(S_{t+1} = s' \mid S_t = s, A_t = a)$

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

(s, s')



$$P_{3,4}^1 = p$$

$$P_{3,1}^2 = 1 - p$$

$$5 \times 5 = 25$$

$$P_{3,5}^2 = p$$

$$P_{3,6}^2 = 0$$

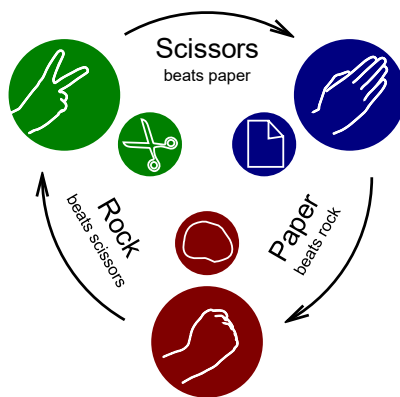
Optimal Policy

• 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

$$\underline{\underline{\pi}}: s \rightarrow \mathcal{A}$$

μ

$$\begin{matrix} \rightarrow 1-\epsilon \\ \rightarrow \underline{\underline{\epsilon}} \end{matrix}$$

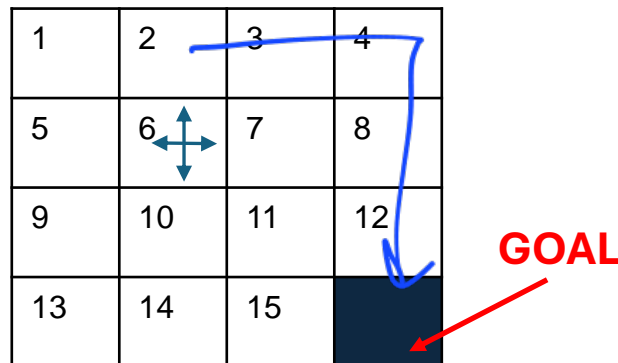
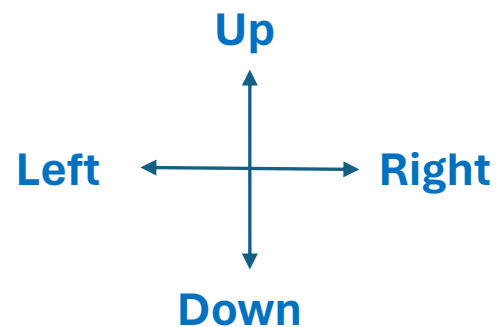
3

- 0 $\rightarrow 1/4$ ✓
- 1 $\rightarrow 1/2$ ✓
- 2 $\rightarrow 1/8$ ✓
- 3 $\rightarrow 1/8$ ✓

explore
explore

How to model your problem as an MDP?

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



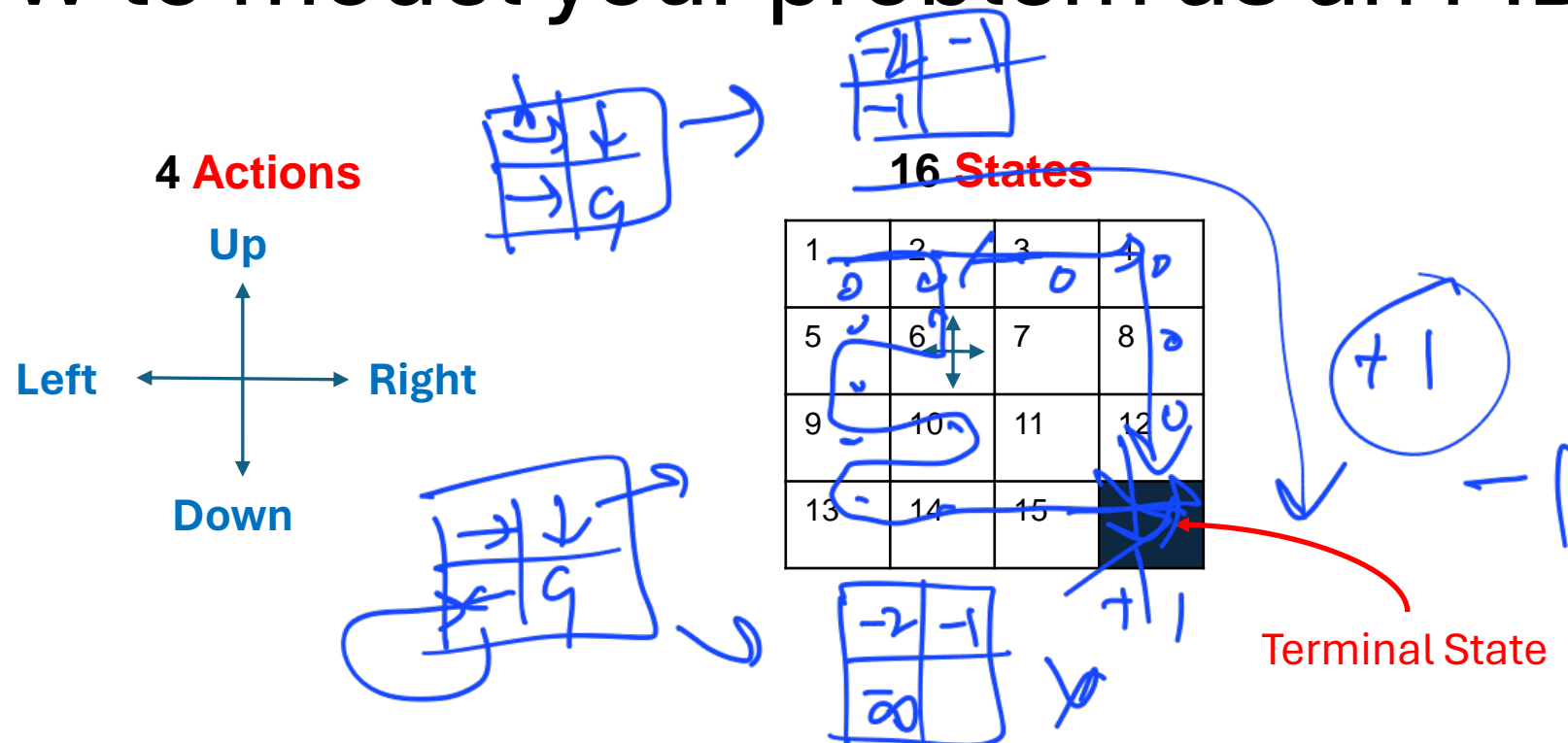
• How to formulate this maze-solving problem as an MDP?

- States ?
- Actions ?
- Rewards ?
- Transition Probabilities ?
- Discount factor ?

Handwritten notes in blue ink:

- S (States)
- A (Actions)
- R_s^a (Rewards)
- $P_{ss'}^a$ (Transition Probabilities)
- γ (Discount factor)

How to model your problem as an MDP?



Rewards

$R_t = -1$ on all transitions

Discount Factor

$\gamma = 1$

$$\pi(s) \rightarrow a$$

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

Verify that optimal policy = shortest path

$$P_{3,4}^R = 1 \quad G_t = R_{t+1} + \gamma V_{t+2} - V_t$$

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?

$$\gamma = 0.9$$

$$0 + \gamma \cdot \underline{1} = 0.9 \checkmark$$

$$\rightarrow s_t = \underline{0} + \gamma \cdot \underline{0} + \gamma^2 \cdot 1 = \underline{0.81} \checkmark \dots$$

$$0 \dots + \underline{1}$$

$$\gamma < 1$$

Bellman Equations

Prof. Subrahmanya Swamy

The diagram illustrates the transformation of a Bellman equation for a state-action pair. On the left, a 2x2 grid represents the action space. The top row contains a right arrow and a down arrow. The bottom row contains a right arrow and a circled 'G'. Above the grid is a circled π . An equals sign follows the grid. To the right is another 2x2 grid representing the matrix form. The top row contains -2 and -1. The bottom row contains -1 and a circled 'G'. Above this grid is a circled V_π .

$$\begin{array}{|c|c|} \hline \rightarrow & \downarrow \\ \hline \rightarrow & \textcircled{G} \\ \hline \end{array} = \begin{array}{|c|c|} \hline -2 & -1 \\ \hline -1 & \textcircled{G} \\ \hline \end{array}$$

Outline

- MDP Dynamics $R_s^a, P_{ss'}^a$
- Policy Dynamics $R_s^\pi, P_{ss'}^\pi$
- Value Function $V_\pi(s)$
- Action-Value Function $Q_\pi(s, a)$
- Bellman Equations

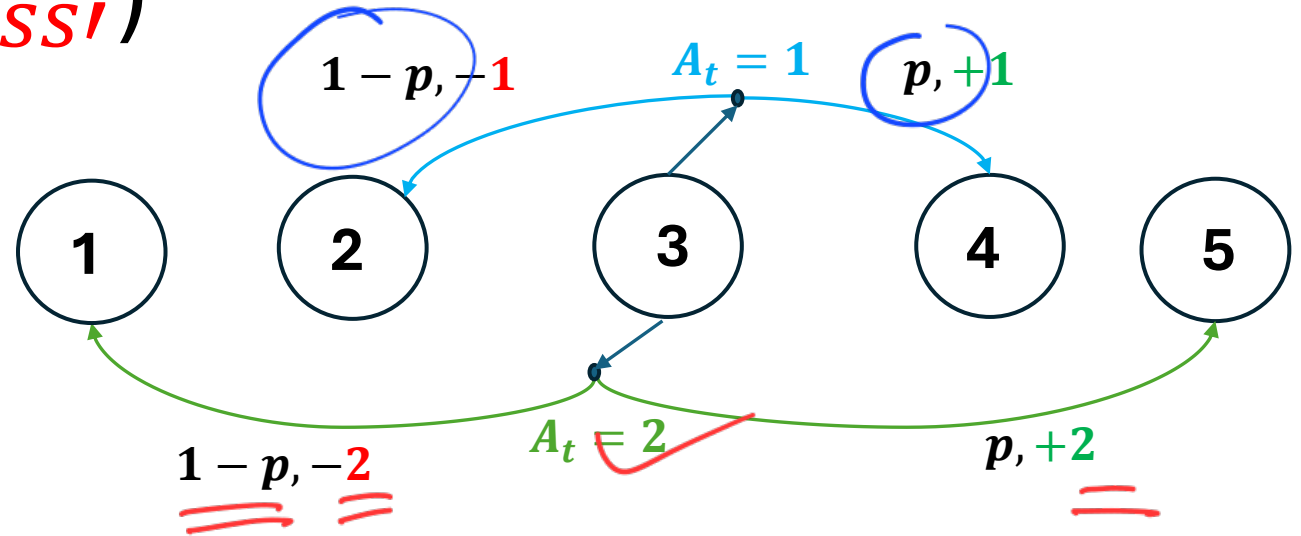


MDP Dynamics ($R_s^a, P_{ss'}^a$)

Transition Probability

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

- Example: $P_{3,5}^2 = p$



Expected Reward

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

- Example:

$$R_3^2 = 2p - 2(1 - p)$$

$$= 4p - 2$$

$$= -1 \quad (\text{if } p = \frac{1}{4})$$

Policy Dynamics $(R_s^\pi, P_{ss'}^\pi)$

Transition Probability ✓

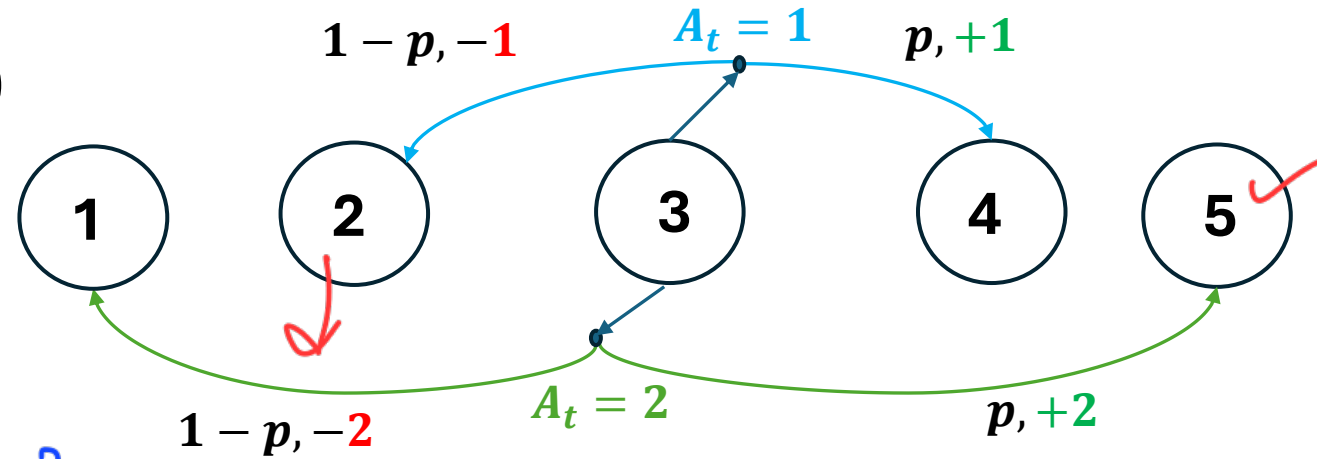
- $P_{ss'}^\pi = \mathbb{P}(S_{t+1} = s' \mid S_t = s, A_t \sim \pi)$
- $= \sum_a \pi(a \mid s) P_{ss'}^a$

$$\frac{1}{4} P_{ss'}^1 + \frac{1}{8} P_{ss'}^2 + \dots$$

Expected Reward

- $R_s^\pi = \mathbb{E}[R_{t+1} \mid S_t = s, A_t \sim \pi]$
- $= \sum_a \pi(a \mid s) R_{ss'}^a$

$$\frac{1}{8} P_{ss'}^2 + \frac{1}{2} P_{ss'}^0$$



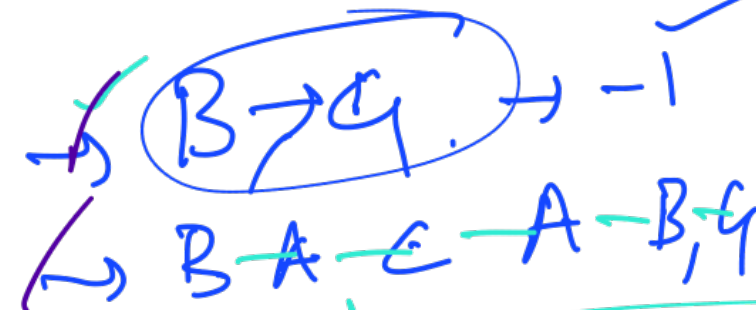
$$\left. \begin{aligned} \pi(a=1 \mid s=2) &= 1/4 \\ \pi(a=2 \mid s=2) &= 1/8 \\ &= 1/8 \\ &= 1/2 \end{aligned} \right\}$$

Value Function ($V_\pi(s)$)

The expected return for following policy π starting from state s

$$V_\pi(s) := \mathbb{E}_\pi[G_t \mid S_t = s]$$

"Bellman eq."



$$\frac{1}{4} \times \frac{1}{4} \times \frac{1}{2} \times \frac{1}{4} \times \frac{1}{4} =$$

$$\sum_{\text{paths}} P(\text{path}) G_t(\text{path})$$

$$E[x] = \sum_x P(x) \cdot x$$

$$\pi(-1) \times \frac{1}{4} \pi(v|B) = \frac{1}{4}$$

$$\frac{1}{4} \times \frac{1}{4}$$

Action-Value Function ($Q_\pi(s, a)$)

The expected return for taking action a in current state s and then following policy π from the next state

$$Q_\pi(s, a) := \mathbb{E}_\pi [G_t \mid S_t = s, A_t = a]$$

Handwritten annotations: A green circle around π in the expectation operator, and a green double underline under the entire expression. A purple arrow points from the state s to a green box labeled 'B', and another purple arrow points from the action a to a green box labeled 'L'.

