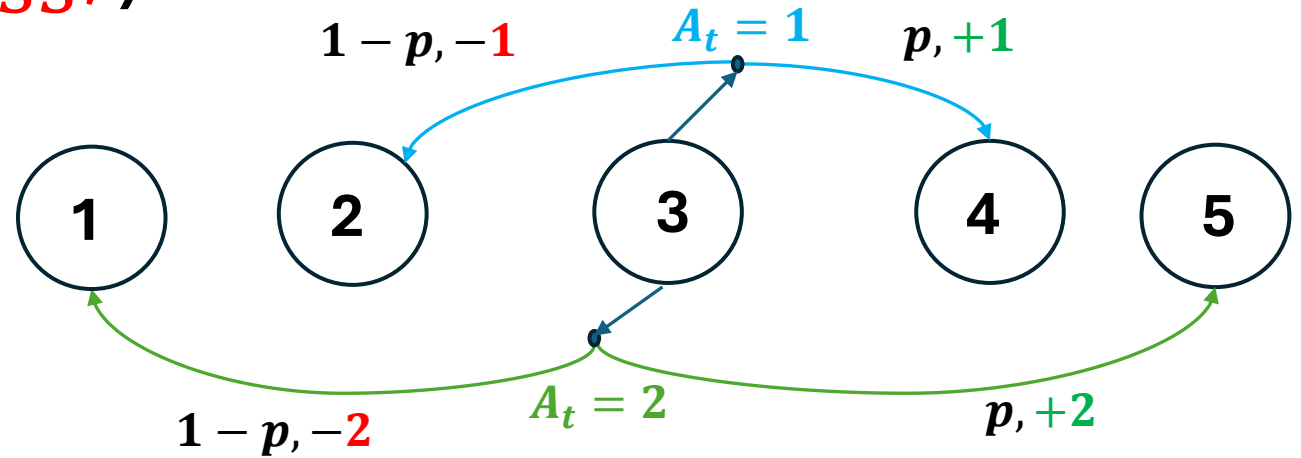# Bellman Equations

Prof. Subrahmanya Swamy

# 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*:

$$R_3^2 = 2\,p - 2\,(1 - p)$$
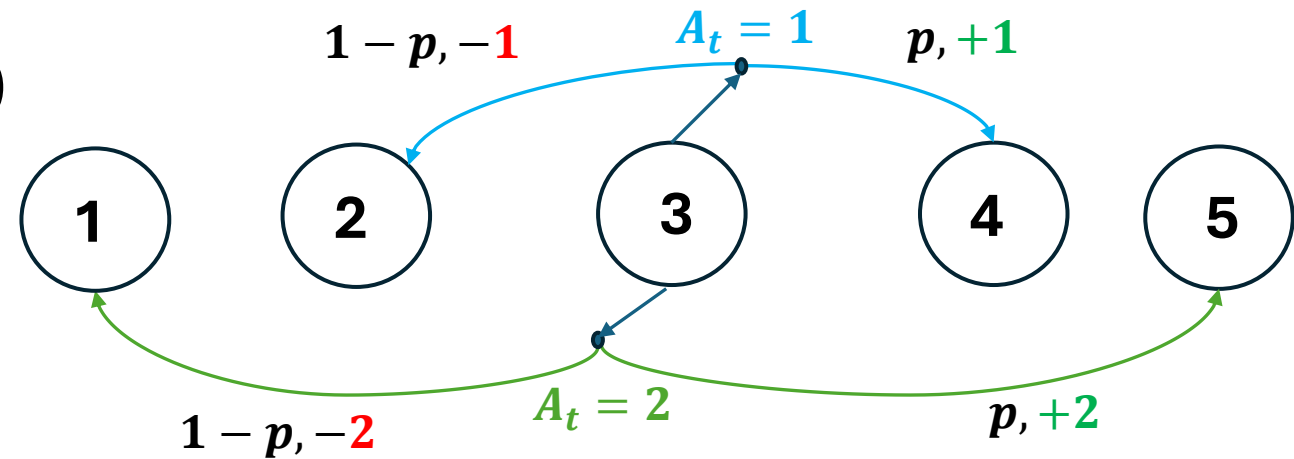$$= 4p - 2$$
$$= -1 \quad (if\ p = \tfrac{1}{4})$$

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

Transition Probability

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



Expected Reward

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

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

The expected return for following policy $\pi$ starting from state $s$

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

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

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

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

# Relating $Q_\pi$ and $V_\pi$

$$V_\pi(s) = \sum_a \pi(a \mid s)\, Q_\pi(s, a)$$

# Relating $Q_\pi$ and $V_\pi$

- $Q_\pi(s, a)$ $= \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a]$
  $= \mathbb{E}_\pi[R_{t+1} + G_{t+1} \mid S_t = s, A_t = a]$
  $= \mathbb{E}_\pi[R_{t+1} \mid S_t = s, A_t = a] + \mathbb{E}_\pi[G_{t+1} \mid S_t = s, A_t = a]$
  $= R_s^a + \sum_{s'} P_{ss'}^a \mathbb{E}_\pi[G_{t+1} \mid S_{t+1} = s', S_t = s, A_t = a]$
  $= R_s^a + \sum_{s'} P_{ss'}^a V_\pi(s')$

- Substitute this in $V_\pi(s) = \sum_a \pi(a \mid s) Q_\pi(s, a)$ to get $V_\pi$ in terms of $V_\pi$

# Bellman Expectation (BE) equation

- $V_\pi$ in terms of $V_\pi$ : (Useful to compute $V_\pi$ from $P_{ss'}^a$ and $R_s^a$)

$$V_\pi(s) = R_s^\pi + \sum_s P_{ss'}^\pi V_\pi(s')$$

**Immediate reward**

**Remaining Return**

▸ $R_s^\pi := \sum_a R_s^a \pi(a \mid s)$

▸ $P_{ss'}^\pi := \sum_a P_{ss'}^a \pi(a \mid s)$