# Bellman Equations

Prof. Subrahmanya Swamy

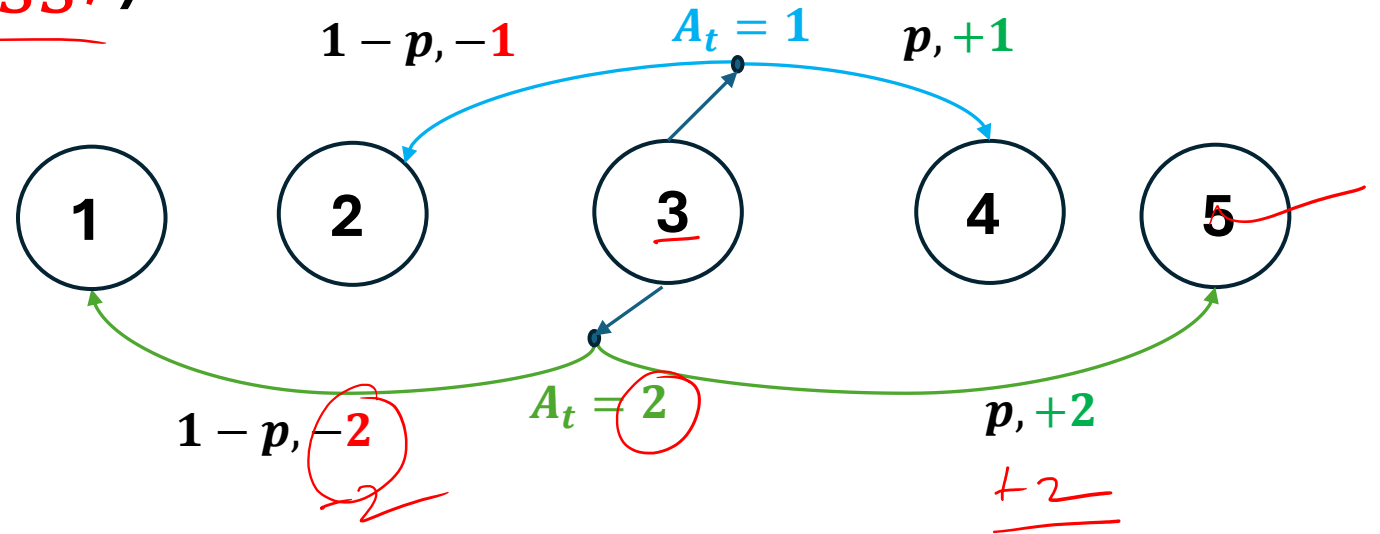# Outline

- MDP Dynamics $\quad R_s^a, P_{ss'}^a$ ✓

- Policy Dynamics $\quad R_s^\pi, P_{ss'}^\pi$ ✓

- Value Function $\quad V_\pi(s)$ ✓

- 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*:

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



$1-p, -1$     $A_t = 1$     $p, +1$

$1-p, -2$     $A_t = 2$     $p, +2$

$+2$

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

### Transition Probability

$\pi(a=2|s=3) = \frac{1}{2}$
$\pi(a=1|s=3) = \frac{1}{2}$

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

$P_{ss'}^\pi = \frac{1}{2} \cdot P_{ss'}^2 + \frac{1}{2} P_{ss'}^1$

$1-p, -1$    $A_t = 1$    $p, +1$

$1-p, -2$    $A_t = 2$    $p, +2$

(1) (2) (3) (4) (5)

### Expected Reward

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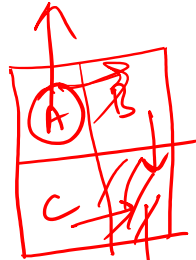- $= \sum_a \pi(a \mid s) \, R_{ss'}^a$

$R^\pi_3$

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

Bellman Expectation Eq.

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

$\pi(\text{Left}|s) = \frac{1}{4}$
$\pi(R|s) = \frac{1}{4}$ ✓
$\pi(D|s) = \frac{1}{4}$
$\pi(U|s) = \frac{1}{4}$

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

$\underbrace{\phantom{V_\pi(s)}}_{V_\pi(A)}$

$\to P(A, \text{Right}, -1, B, \text{Down}, -1, G) \to \boxed{-2}$

$\to A, \text{UP}, -1, A, \text{Down}, -1, C, \text{Right}, -1, G \to \boxed{-3}$

↑↑↑↑↑
A A A A A ... B G

$\pi(\text{Right}) \times \pi(\text{Down}) =$

$\frac{1}{4} \times \frac{1}{4} = \frac{1}{16}$

$-2 \times \frac{1}{16} \quad -3 \times \left(\frac{1}{4}\right)^3 -$

# 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

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

$$V_\pi(s) \qquad Q_\pi(s,a)$$

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

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

$$\pi \rightarrow \begin{matrix} a_1 \rightarrow \pi(a_1 \mid s) \\ a_2 \rightarrow \pi(a_2 \mid s) \end{matrix}$$

$$V_\pi(s) = \pi(a_1 \mid s) Q_\pi(s, a_1) + \pi(a_2 \mid s) Q_\pi(s, a_2)$$

$$Q_\pi(s,a) \qquad V_\pi(s) \qquad G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots$$

$$G_{t+1} = R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} \cdots$$

$$Q_\pi(s,a) = E_\pi\left[ G_t \mid S_t = s, A_t = a \right]$$

$$E[X+Y] = E[X] + E[Y]$$

$$= E_\pi\left[ R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a \right]$$

$$= E_\pi\left[ R_{t+1} \mid S_t = s, A_t = a \right] + \gamma E_\pi\left[ G_{t+1} \mid S_t = s, A_t = a \right]$$

$$R_s^a + \gamma$$

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

$V_\pi \leftrightarrow Q_\pi$

$V_\pi \leftrightarrow$

| A | B |
|---|---|
| C | |

$V_\pi(A)$
$V_\pi(B)$
$V_\pi(C)$

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

$\quad = \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a]$

$\quad = \mathbb{E}_\pi[R_{t+1} \mid S_t = s, A_t = a] + \mathbb{E}_\pi[G_{t+1} \mid S_t = s, A_t = a]$

$\quad = R_s^a + \gamma \sum_{s'} P_{ss'}^a \mathbb{E}_\pi[G_{t+1} \mid S_{t+1} = s', S_t = s, A_t = a]$

$Q_\pi(s,a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a V_\pi(s') \quad —①$

Markov property

- Substitute this in $V_\pi(s) = \sum_a \pi(a \mid s) Q_\pi(s,a)$ to get $V_\pi$ in terms of $V_\pi$
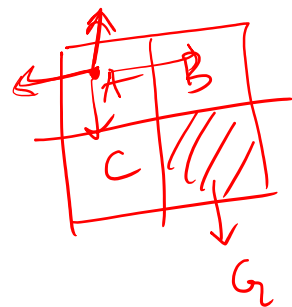
②

# 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**

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

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

Example

$\pi \to$ Uniform Random
Policy



$G$

$BE$ $\boxed{V_\pi(s)} = R_s^\pi + \gamma \underset{s'}{\boxed{\sum}} P_{ss'}^\pi \underline{V_\pi(s')}$

① ②

$\uparrow$ Immediate Reward

$\downarrow$ Remaining Return

3 unknowns
3 linear equations

$V_\pi(A)$

$10^6 \times 10^6$

① $R_A^\pi = -1$

$P_{AB}^\pi = \frac{1}{4}$, $P_{AC}^\pi = \frac{1}{4}$, $P_{A,A}^\pi = \frac{1}{2}$

$\boxed{Ax = b}$
$\boxed{x = A^{-1}b}$

②

① $\begin{cases} V_\pi(A) = -1 + \gamma \left( P_{AB}^\pi V_\pi(B) + P_{AC}^\pi V_\pi(C) + P_{AA}^\pi V_\pi(A) \right) \\ \\ V_\pi(B) = -1 + \cdots \\ \\ V_\pi(C) = -1 + \cdots \end{cases}$