

Bellman Expectation Bellman Optimality Iterative Policy Evaluation

Prof. Subrahmanya Swamy

Bellman Expectation Equations: Numerical Example

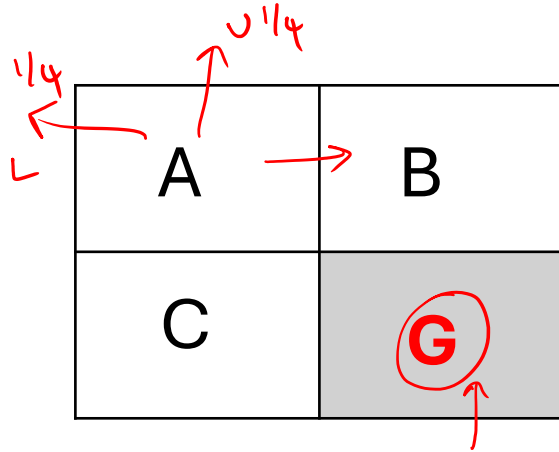
Bellman Expectation (BE) equation

$$\checkmark \quad \underline{V}_{\pi}(\underline{s}) = \underline{R}_{\underline{s}}^{\pi} + \sum_{\underline{s}'} \underbrace{P_{\underline{s}\underline{s}'}^{\pi}}_{\text{Remaining Return}} \underline{V}_{\pi}(\underline{s}')$$

↑ Immediate reward
↑

To find the value function \underline{V}_{π} of a given policy π

Grid Example



- Deterministic state transitions $P_{A,B}^R = 1$
- $R_t = -1$ on all transitions
- Terminal state value $V_\pi(G) = 0$
- Discount factor $\gamma = 1$
- "Uniform" Random Policy π
 - U $\rightarrow 1/4$
 - R $\rightarrow 1/4$
 - L $\rightarrow 1/4$

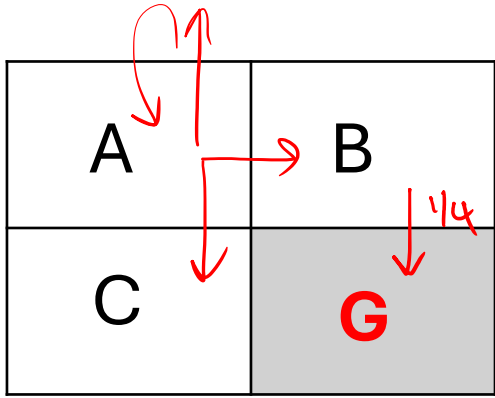
Policy Dynamics:

$$\underline{P_{A,A}^\pi} = \frac{1}{2}, \quad \underline{P_{A,B}^\pi} = \frac{1}{4}, \quad P_{A,C}^\pi = \frac{1}{4}$$

$$P_{B,A}^\pi = \frac{1}{4}, \quad P_{B,B}^\pi = \frac{1}{2}, \quad P_{B,G}^\pi = \frac{1}{4}$$

$$P_{C,A}^\pi = \frac{1}{4}, \quad P_{C,G}^\pi = \frac{1}{4}, \quad P_{C,C}^\pi = \frac{1}{2}$$

Bellman Expectation



BE Eqn.

$$V_{\pi}(s) = \underset{-1}{R_{\underset{\downarrow}{s}}^{\pi}} + \sum_{\underset{\downarrow}{s'}} P_{ss'}^{\pi} V_{\pi}(s')$$

A: $V_{\pi}(\underline{A}) = \underline{-1} + \frac{1}{4}V_{\pi}(\underline{B}) + \frac{1}{4}V_{\pi}(\underline{C}) + \overset{P_{AA}^{\pi}}{\frac{1}{2}}V_{\pi}(\underline{A})$ ✓

B: $V_{\pi}(\underline{B}) = -1 + \frac{1}{4}V_{\pi}(\underline{A}) + \frac{1}{4}\underline{V_{\pi}(\underline{G})} + \frac{1}{2}V_{\pi}(\underline{B})$

C: $V_{\pi}(\underline{C}) = -1 + \frac{1}{4}V_{\pi}(\underline{A}) + \frac{1}{4}V_{\pi}(\underline{G}) + \frac{1}{2}V_{\pi}(\underline{C})$

Matrix form

$$\mathbf{A:} \quad V_{\pi}(A) = -1 + \frac{1}{4}V_{\pi}(B) + \frac{1}{4}V_{\pi}(C) + \frac{1}{2}V_{\pi}(A)$$

$$\mathbf{B:} \quad V_{\pi}(B) = -1 + \frac{1}{4}V_{\pi}(A) + \frac{1}{4}V_{\pi}(\underline{G}) + \frac{1}{2}V_{\pi}(B)$$

$$\mathbf{C:} \quad V_{\pi}(C) = -1 + \frac{1}{4}V_{\pi}(A) + \frac{1}{4}V_{\pi}(\underline{G}) + \frac{1}{2}V_{\pi}(C)$$

$$\Rightarrow \begin{bmatrix} -\frac{1}{2} & \frac{1}{4} & \frac{1}{4} \\ \frac{1}{4} & -\frac{1}{2} & 0 \\ \frac{1}{4} & 0 & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} V_{\pi}(A) \\ V_{\pi}(B) \\ V_{\pi}(C) \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

Matrix form

Solving the matrix equation gives us

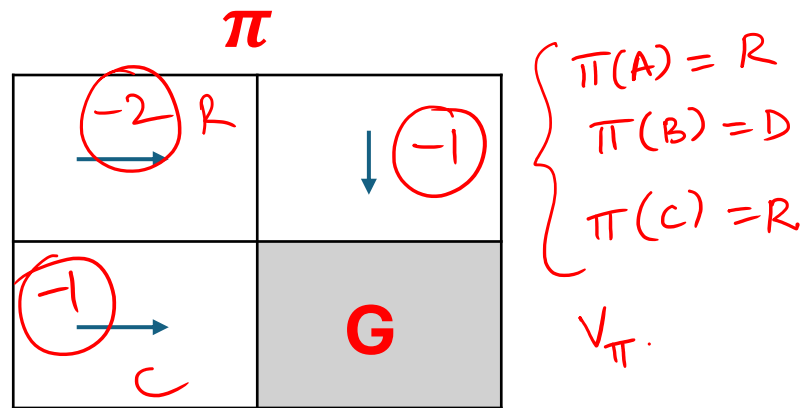
V_{π}

-8 ^A	-6 ^B
-6 ^C	0

Uniform random Policy

Exercise

- Use BE equations and compute the value function for the policy shown in the figure



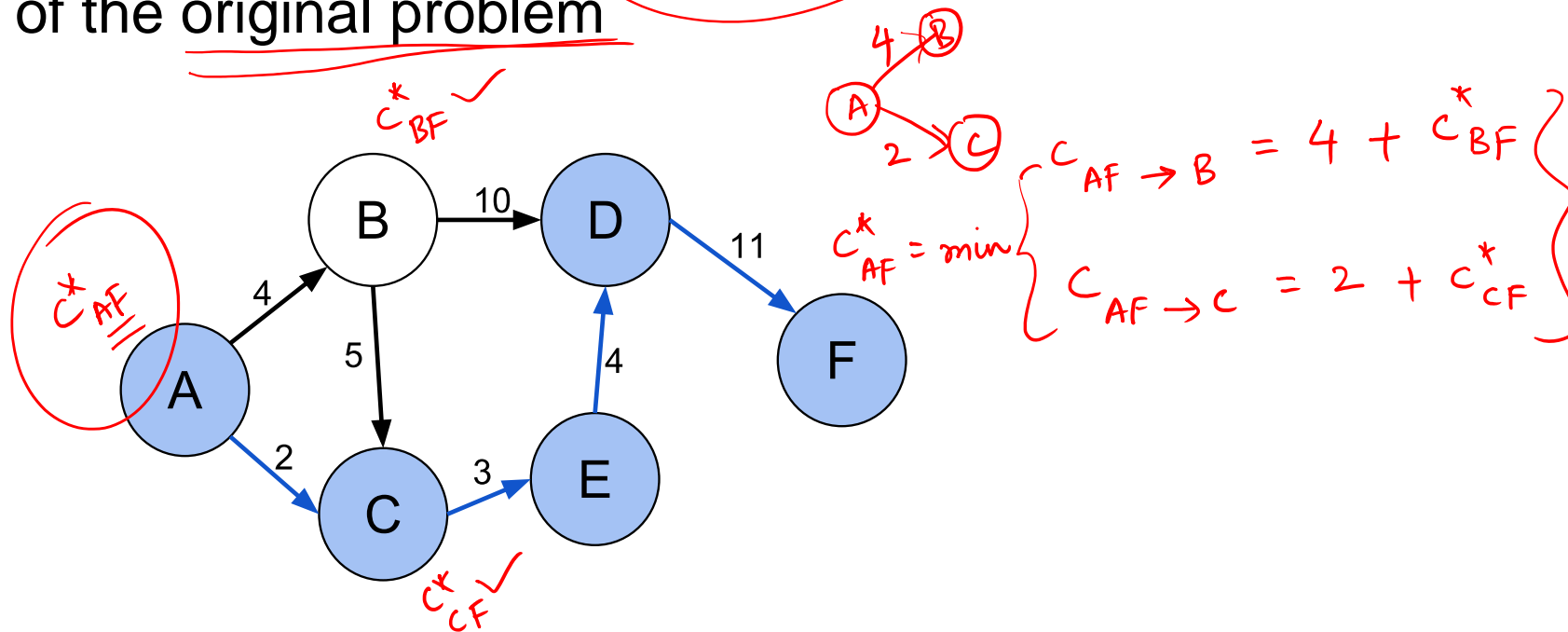
Bellman Optimality Equations

Bellman Optimality Equations

- **Bellman Expectation :** To find V_π for a given policy π <sup>↓
unif</sup>
- **Bellman Optimality :** To find optimal policy π^*

Optimal substructure

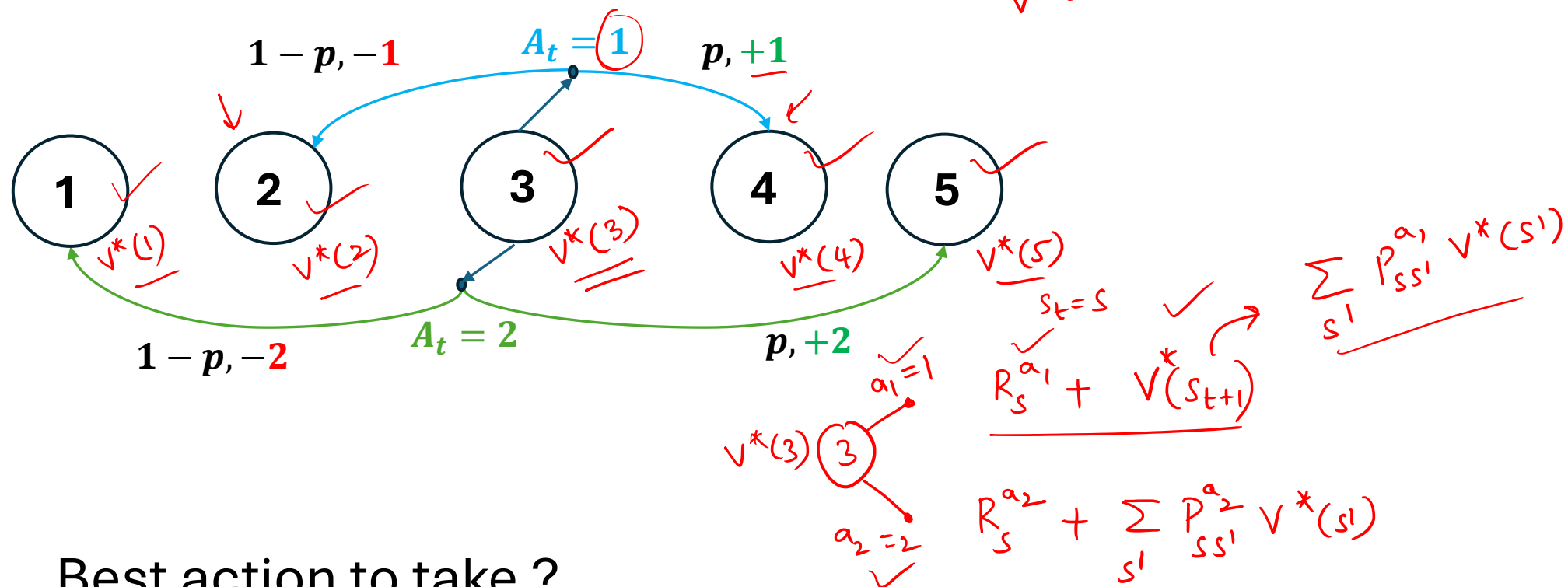
Optimal solutions of subproblems can be used to find the optimal solution of the original problem



The shortest cost for A -> F

Can be **found from** the shortest costs of **B -> C**, **C -> F**

Optimal Substructure in MDP



Best action to take ?

Bellman Optimality (BO) equation

$$Q_{\pi}(s_t, a) = R_s^a + \sum_{s'} P_{ss'}^a V^{\pi}(s')$$

$$V^*(s) = \max_a \left[\underbrace{R_s^a}_{\text{Reward for action } a} + \sum_{s'} P_{ss'}^a \underbrace{V^*(s')}_{\text{Remaining Reward from next state}} \right] = Q_{\pi^*}(s, a)$$

Using the definition of Q-function

We can equivalently write it as

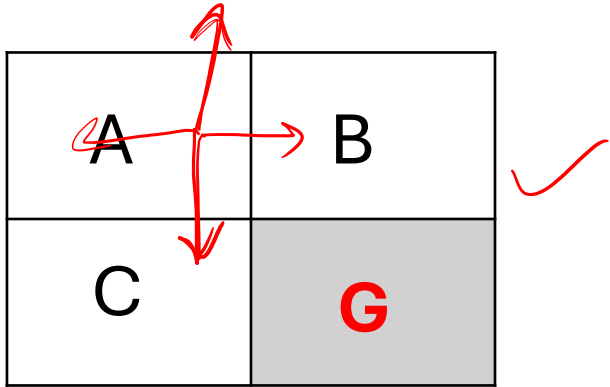
$$V^*(s) = \max_a Q^*(s, a)$$

Optimal Policy from V^*



$$\underline{\pi^*(s)} = \arg \max_a R_s^a + \sum_{s'} P_{ss'}^a V^*(s')$$

Example: Verify BO equations



$V_*(s)$

-2 _A	-1 _B
-1 _C	0

$$\rightarrow \underline{V^*(s)} = \max_a \left\{ \underline{R_s^a} + \sum_{s'} \underline{P_{ss'}^a} \underline{V^*(s')} \right\}^{BO}$$

A:

$$\begin{array}{lcl}
 \begin{array}{c} V^*(A) \\ \parallel \\ -2 \end{array} & \begin{array}{l} \xrightarrow{UP} \\ \xrightarrow{D} \\ \xrightarrow{L} \\ \xrightarrow{R} \end{array} & \begin{array}{l} -1 + V^*(A) \cdot 1 = -1 - 2 = -3 \checkmark \\ -1 + V^*(C) = -1 - 1 = -2 \checkmark \\ -1 + V^*(A) = -1 - 2 = -3 \checkmark \\ -1 + V^*(B) = -1 + -1 = -2 \checkmark \end{array} \\
 & & = \max\{-3, -2, -3, -2\} = \underline{-2}
 \end{array}$$

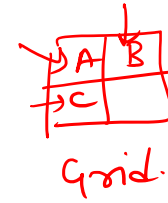
Iterative Policy Evaluation

Iterative Policy Evaluation

- For large state spaces

V^π

Bellman Expectation Eq.



- Policy Evaluation: Iteratively apply BE equation $V_{k+1}(s) = R_s^\pi + \sum_{s'} P_{ss'}^\pi V_k(s')$

V^π

$\pi \rightarrow$ uniform

$$V_0 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad V_{k+1}(s) = R_s^\pi + \sum_{s'} P_{ss'}^\pi V_k(s') \quad \checkmark$$

$$V_1(A) = -1 + P_{A,A}^\pi V_0(A) + P_{A,B}^\pi V_0(B) + P_{A,C}^\pi V_0(C)$$

$$V_1(A) = -1 + \frac{1}{2}(0) + \frac{1}{4}(0) + \frac{1}{4}(0)$$

$$V_0 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$



$$V_0 \xrightarrow{V_1} \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} \xrightarrow{V_2} \begin{bmatrix} \\ \\ \end{bmatrix} \xrightarrow{V_3} \dots V_{k+1}(s) = R_s^\pi + \sum_{s'} P_{ss'}^\pi V_k(s')$$

$$\dots V_{k+1} = V_k.$$

$$\Downarrow$$

$$V_k = \underline{V_\pi}$$

$$V_1(A) = -1 + 0 + 0 + 0$$

$$V_1(B) = -1 + 0 + 0 + 0$$

$$V_1(C) = -1$$

$k=2$

$$V_2(A) = -1 + \frac{1}{2} V_1(A) + \frac{1}{4} V_1(B) + \frac{1}{4} V_1(C)$$

$$= -1 + \frac{1}{2}(-1) + \frac{1}{4}(-1) + \frac{1}{4}(-1)$$

$$= -2$$

$$V_2(B) = -1 + \frac{1}{4} V_1(C) + \frac{1}{4} V_1(A) + \frac{1}{2} V_1(B)$$

$$= -1 + 0 - \frac{1}{4} - \frac{1}{2}$$

$$V_2(C) =$$