

Bellman Expectation Bellman Optimality Iterative Policy Evaluation

Prof. Subrahmanya Swamy

Bellman Expectation Equations: Numerical Example

Bellman Expectation (BE) equation

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

Immediate reward Remaining Return

To find the value function V_{π} of a given policy π

Grid Example

A	B
C	G

- **Deterministic** state transitions
- $R_t = -1$ on all transitions
- Terminal state value $V_\pi(G) = 0$
- Discount factor $\gamma = 1$
- **Uniform** Random Policy π

Policy Dynamics:

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

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

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

Bellman Expectation

A	B
C	G

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

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

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

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

Matrix form

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

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

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

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

Matrix form

Solving the matrix equation gives us

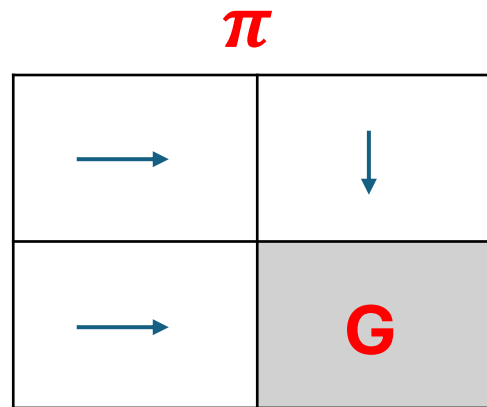
V_{π}

-8	-6
-6	0

Uniform random Policy

Exercise

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



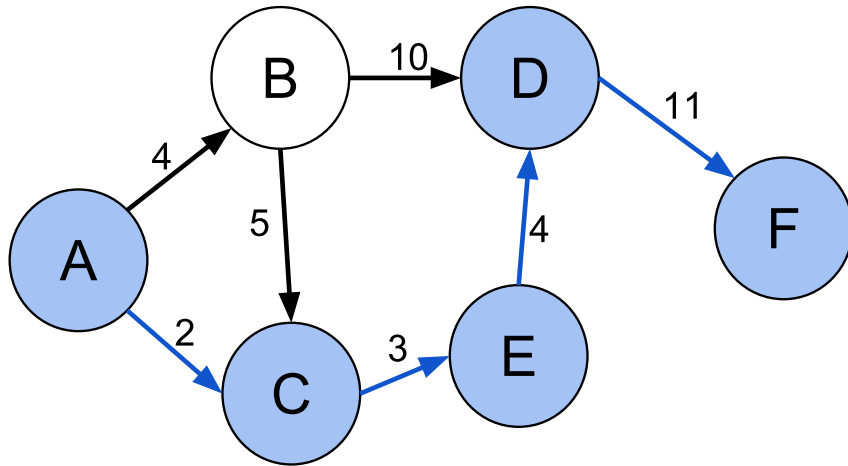
Bellman Optimality Equations

Bellman Optimality Equations

- **Bellman Expectation :** To find V_π for a **given** policy π
- **Bellman Optimality :** To find optimal policy π^*

Optimal substructure

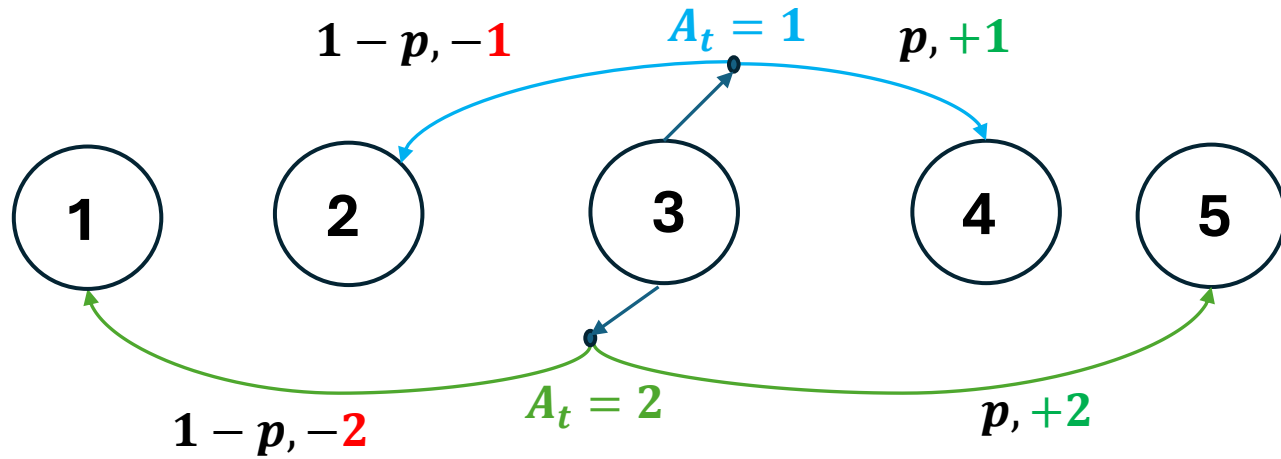
Optimal solutions of **subproblems** can be used to find the optimal solution of the original problem



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

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

Reward for
action a

Remaining
Reward from
next state

Using the definition of Q-function

We can equivalently write it as

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

Optimal Policy from V^*

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

Example: Verify BO equations

A	B
C	G

V_*

-2	-1
-1	0

Iterative Policy Evaluation

Iterative Policy Evaluation

- ▶ Large state spaces:
 - ▶ *Issue:* Solving Bellman expectation equations using matrix inversion is intractable
 - ▶ *Solution:* Use iterative policy evaluation
- ▶ Iterative Policy Evaluation: Iteratively apply BE equation

$$V_{k+1}(s) = R_s^\pi + \sum_{s'} P_{ss'}^\pi V_k(s')$$