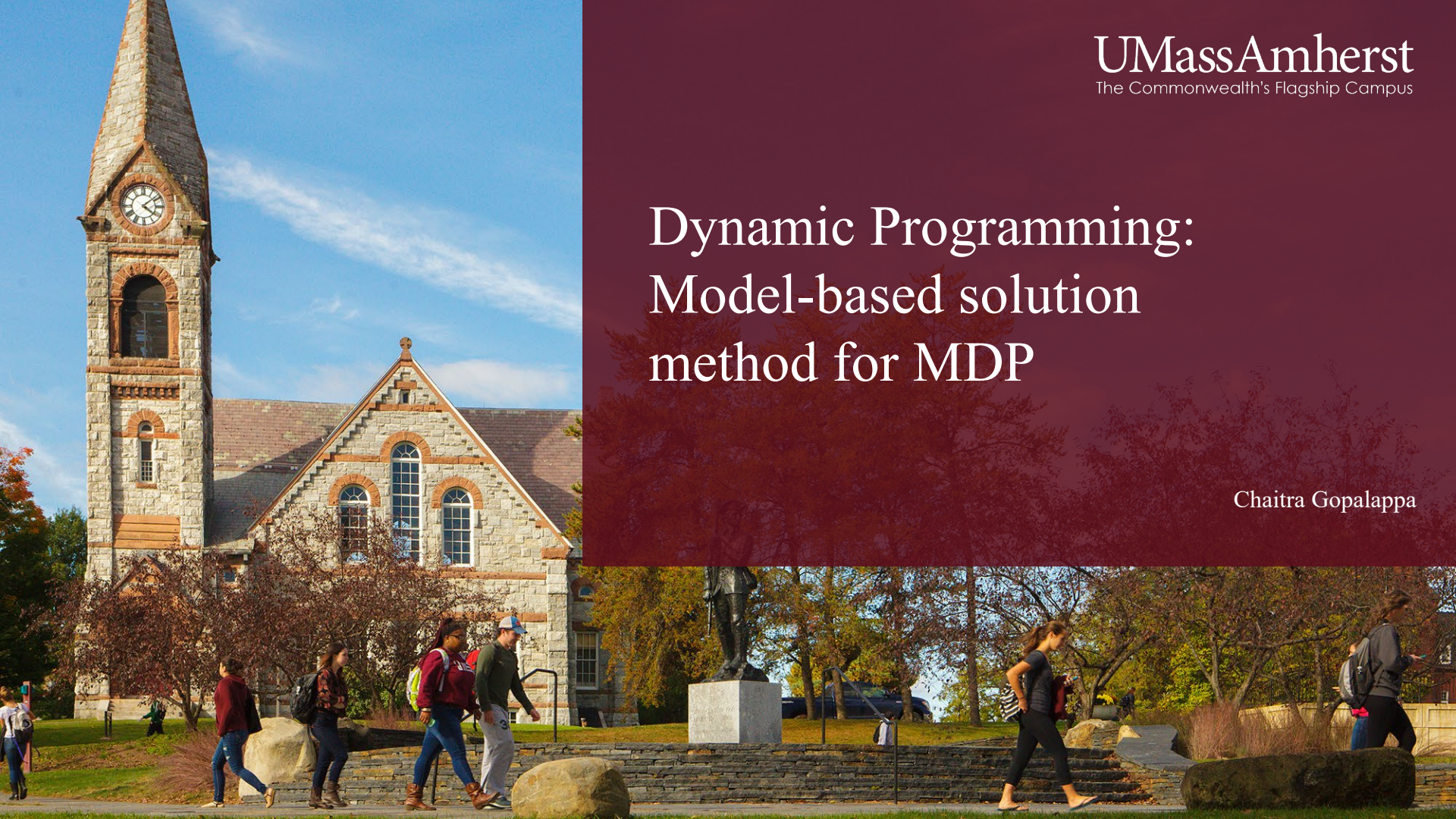


Dynamic Programming: Model-based solution method for MDP

Chaitra Gopalappa



Reviewing MDP terminology for a refresher

Caution: Unless otherwise specified, here on:

- π
 - is a policy (not steady state distribution like in Markov chain)
- Value function:
 - $v_{\pi}(s) = \mathbb{E}_{\pi}[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s]; \forall s \in S$

Objective of MDP: Find policy with maximum state-value function

- $v_*(s) = \max_{\pi} v_{\pi}(s), \forall s \in \mathcal{S}$
 - $v_{\pi}(s)$: value function of state s under policy π
 - It is the **expected** reward when **starting** in state s and following some policy π
 - $v_*(s)$: value function of state s under optimal policy π^*
- Later-on we will use other objective functions

Expand state-value function

$$\begin{aligned} v_{\pi}(s) &= \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right] \\ &= \mathbb{E}_{\pi} [R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \cdots | S_t = s, \pi] \\ &= \mathbb{E}_{\pi} [R_t + \gamma (R_{t+1} + \gamma R_{t+2} + \cdots) | S_t = s, \pi] \\ &= \mathbb{E}_{\pi} [R_t + \gamma (v_{\pi}(S_{t+1})) | S_t = s, \pi] \\ &= \sum_{s'} p(s, \pi(s), s') [r(s, \pi(s), s') + \gamma v_{\pi}(s')] \end{aligned}$$

Dynamic programming

(Chapter 4, Sutton and Barto)

DP Algorithms

- Policy iteration solved as system of equations
- Policy iteration (this is the one generally referred to)
- Value iteration

Bellman equation for a “fixed” (deterministic) policy π

$$V(s) = \sum_{s'} p(s, \pi(s), s') [r + \gamma V(s')]$$

- $r = r(s, \pi(s), s')$; r here is deterministic
- Note: if there is randomness in value of r we rewrite the value function as
 - $V(s) = \sum_{s', r} [p(s, \pi(s), s', r) [r + \gamma V(s')]]$

Bellman “optimality” equation

- Let v^* be the optimal value function for the MDP. The function v^* satisfies, for each $s \in S$, the following
 - $v^*(s) = \max_{a \in A} \sum_{s', r} [p(s', r | s, a) \cdot (r(s, a, s') + v^*(s'))]$
 - Or if there is no randomness in rewards,
 - $v^*(s) = \max_{a \in A} \sum_{s'} [p(s' | s, a) \cdot (r(s, a, s') + v^*(s'))]$
- Furthermore, it is the only function satisfying the property

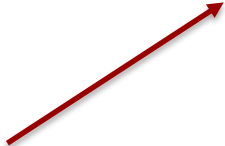
Some basic theories

- Definition: A **stationary policy**, defined by the function π takes actions $\pi(i)$ at time n if $X_n = i$, independent of previous states, previous actions, and time n
- Property: if the state space S is finite, there exists a stationary policy that solves the problem
 - $v^*(s) = \max_{\pi} v_{\pi}(s), \forall s \in S$
- Definition: A function is called invariant with respect to an operation, if the operation does not vary the function.
 - If the invariant function is unique, then it is called a 'fixed point' for the operation

Some basic theories

- Definition: A **stationary policy**, defined by the function π takes actions $\pi(i)$ at time n if $X_n = i$ independent of previous states, previous actions, and time n
- Property: if the state space S is finite, there exists a stationary policy that solves the problem
 - $v^*(s) = \max_{\pi} v_{\pi}(s), \forall s \in S$
- Definition: A function is called invariant with respect to an operation, if the operation does not vary the function.
 - If the invariant function is unique, then it is called a 'fixed point' for the operation

Note: similarity in steady state/
stationary distribution (ρ) of a Markov process,
 $\rho P = \rho$
Multiplying by P does not change the value of ρ , so ρ is the fixed point for the MC defined by P



Fixed point theorem to solve for Bellman operators

- Solving for Bellman equation for a “**given**” (deterministic) policy π

$$V(s) = \sum_{s'} p(s, \pi(s), s') [r + \gamma V(s')]$$

Equivalent to writing: $V_\pi = B_\pi V_\pi$

B_π : Bellman operator for a policy π ;

Applying Banach fixed-point theorem, for a given policy π , we can solve for V_π by starting at some random value V

- Solving for Bellman **optimality** equation

$$v^*(s) = \max_{a \in A} \sum_{s'} [p(s' | s, a) \cdot (r(s, a, s') + \gamma v^*(s'))]$$

Equivalent to writing: $V_{\pi^*}^* = B_{\pi^*}^* V_{\pi^*}^*$ or $V^* = B^* V^*$

B^* : Bellman optimality operator (corresponding to optimal policy π^* ;

Applying Banach fixed-point theorem we can solve for $v^*(s)$ by starting at some random value V

Policy iteration solved as system of equations

1. Initialize $\pi(s) \in \mathcal{A}(s)$; arbitrarily $\forall s \in \mathcal{S}$
2. Policy evaluation
 - *Solve for $V(s)$, each $s \in \mathcal{S}$, by solving following system of equations*
 - $V(s) = \sum_{s',r} [p(s',r | s, \pi(s)) \cdot (r + \gamma \cdot V(s'))]$; $r = R(s, \pi(s), s')$
 - 3. Policy improvement
 - *For each $s \in \mathcal{S}$:*
 - $\pi'(s) = \operatorname{argmax}_a \sum_{s',r} [p(s',r | s, a) \cdot (r + \gamma \cdot V(s'))]$
 - If $\pi'(s) = \pi(s)$, stop and return $v_* = V(s), \pi_* = \pi$, else goto step 2

Policy Iteration DP
 $\{X_t, D_t\}_{t=0}^n$ is a MDP represented by a tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$; $\mathcal{S} = \{1, 2\}$
 $A = \{a_1, a_2\}$
 $P_{a_1} = \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix}$; $P_{a_2} = \begin{bmatrix} 0.9 & 0.1 \\ 0.2 & 0.8 \end{bmatrix}$; $R_{a_1} = \begin{bmatrix} 6 & -5 \\ 7 & 12 \end{bmatrix}$; $R_{a_2} = \begin{bmatrix} 10 & 17 \\ -14 & 13 \end{bmatrix}$; Find optimal policy

- ① Let $\pi = [a_1, a_1]$; initializing π ; $\lambda =$ discounting factor
 ② Policy evaluation $V(s) = \sum_{s', r} [P(s', r | s, \pi(s)) \cdot (r + \gamma V(s'))]$; $r = R(s, \pi(s), s')$
 $V(1) = [0.7 \ 0.3] \begin{bmatrix} 6 + \gamma V(1) \\ -5 + \gamma V(2) \end{bmatrix} = 2.7 + [0.7V(1) + 0.3V(2)]\lambda$
 $V(2) = [0.4 \ 0.6] \begin{bmatrix} 7 + \gamma V(1) \\ 12 + \gamma V(2) \end{bmatrix} = 10 + [0.4V(1) + 0.6V(2)]\lambda$
 Solve for $V(1)$ and $V(2) \Rightarrow V(1) = 54$; $V(2) = 64$

- ③ Policy improvement
 $\pi'(s) = \arg \max_a \sum_{s', r} [P(s', r | s, a) \cdot (r + \gamma V(s'))]$
 $\pi'(1) = \arg \max_a \left\{ [0.7 \ 0.3] \begin{bmatrix} 6 + \gamma 54 \\ -5 + \gamma 64 \end{bmatrix}; [0.9 \ 0.1] \begin{bmatrix} 10 + \gamma 54 \\ 17 + \gamma 64 \end{bmatrix} \right\}$
 $= \arg \max_a \left\{ 54; 60.2 \right\}$
 $= a_2$
 $\pi'(2) = \arg \max_a \left\{ [0.4 \ 0.6] \begin{bmatrix} 7 + \gamma 54 \\ 12 + \gamma 64 \end{bmatrix}; [0.2 \ 0.8] \begin{bmatrix} -14 + \gamma 54 \\ 13 + \gamma 64 \end{bmatrix} \right\}$
 $= \arg \max_a \left\{ 64; 63.4 \right\}$

$\pi' = [a_2, a_1] \Rightarrow \pi \neq \pi' \Rightarrow$ Set $\pi = \pi'$ goto step 2

- ② policy evaluation: $\pi = [a_2, a_1]$
 $V(1) = [0.9 \ 0.1] \begin{bmatrix} 10 + \gamma V(1) \\ 17 + \gamma V(2) \end{bmatrix} = 10.7 + \gamma [0.9V(1) + 0.1V(2)]$
 $V(2) = [0.4 \ 0.6] \begin{bmatrix} 7 + \gamma V(1) \\ 12 + \gamma V(2) \end{bmatrix} = 10 + \gamma [0.4V(1) + 0.6V(2)]$
 Two equations, two unknowns
 Solve for $V(1)$ and $V(2)$; $V(1) = 105.85$; $V(2) = 104.58$

- ③ policy improvement
 $\pi'(1) = \arg \max_a \left\{ [0.7 \ 0.3] \begin{bmatrix} 6 + \gamma 105.85 \\ -5 + \gamma 104.58 \end{bmatrix}; [0.9 \ 0.1] \begin{bmatrix} 10 + \gamma 105.85 \\ 17 + \gamma 104.58 \end{bmatrix} \right\}$
 $= \arg \max_a \left\{ 97.6; 105.85 \right\}$
 $= a_2$
 $\pi'(2) = \arg \max_a \left\{ [0.4 \ 0.6] \begin{bmatrix} 7 + \gamma 105.85 \\ 12 + \gamma 104.58 \end{bmatrix}; [0.2 \ 0.8] \begin{bmatrix} -14 + \gamma 105.85 \\ 13 + \gamma 104.58 \end{bmatrix} \right\}$
 $= \arg \max_a \left\{ 104.58; 101.95 \right\}$
 $= a_1$
 $\pi' = [a_2, a_1] \Rightarrow \pi = \pi' \Rightarrow \pi = [a_2, a_1]$ is the optimal policy

Fixed point theorem to solve for Bellman operators

- Solving for Bellman equation for a “**given**” (deterministic) policy π

$$V(s) = \sum_{s'} p(s, \pi(s), s') [r + \gamma V(s')]$$

Equivaent to writing: $V_\pi = B_\pi V_\pi$

B_π : Bellman operator for a policy π ;

Applying Banach fixed-point theorem, for a given policy π , we can solve for V_π by starting at some random value V

1. Initialize $V(s) \in \mathbb{R}, \pi(s) \in \mathcal{A}(s)$; arbitrarily $\forall s \in \mathcal{S}$; set tolerance $\theta (=1e-6)$

2. Policy evaluation

- *Loop while $\Delta > \theta$:*
 - $\Delta \leftarrow 0$
 - *Loop for each $s \in \mathcal{S}$*
 - $v \leftarrow V(s)$
 - $V(s) = \sum_{s',r} p(s', r | s, \pi(s)) [r + \gamma V(s')]$
 - $\Delta \leftarrow \max(\Delta, |v - V(s)|)$

Applies fixed-point theorem
to solve Bellman equation for
a *fixed* policy
Note: This can be applied to
any policy

• 3. Policy improvement

- *For each $s \in \mathcal{S}$:*
 - $\pi'(s) = \operatorname{argmax}_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$
- If $\pi'(s) = \pi(s)$, stop and return $v_* = V(s), \pi_* = \pi$, else goto step 2

Finite horizon v infinite time horizon

- Infinite time horizon:
 - e.g., What is optimal inventory ordering policy (DP assignment)
 - Discounting bounds the value function; no changes needed in formulation
 - What if discounting is not preferred? (discounting gives less value to future rewards)
 - $V(s) = \sum_{s',r} [p(s',r | s, \pi(s)) \cdot (R(s, \pi(s), s') + \gamma \cdot V(s'))]$ – \bar{V} ; \bar{V} = average reward associated with the policy under evaluation (unknown); 1 more unknown than number of equation \rightarrow set one of $V(s) = 0$ and solve for the other elements
- Finite time horizon:
 - e.g., game (win, lose); health progression over a lifetime; robotic tasks
 - Make the final state a terminal state ($\Pr(\text{terminal state}, \text{terminal state}) = 1$)
 - Make reward for $R(\text{terminal state}, a, \text{terminal state}) = 0$; you can give one time-reward for the transition, e.g., $R(:, a, \text{win}_{\text{state}}) = \text{high reward}$
 - You may choose to make discounting factor = 1; as time is sufficiently small, it bounds the value function

Fixed-point theorem

- In **policy iteration** we first did policy evaluation
 - Based on Bellman's equation, for any policy π , there exist a fixed-point solution V
 - $V(s) = \sum_{s',r} [p(s', r | s, \pi(s)) \cdot (R(s, \pi(s), s') + \gamma \cdot V(s'))]$
 - Therefore, for any π , we can solve for $V(s)$, by starting with some arbitrary value and iterating through until we arrive at the fixed point (solution converges)
- Applying the fixed-point theorem, we also have
 - If v_* is the optimal value function for a given MDP, then it satisfies the following equation
 - $v_*(s) = \max_{a \in \mathcal{A}} \sum_{s',r} [p(s', r | s, a) \cdot (R(s, a, s') + \gamma \cdot v_*(s'))]$
 - Then why not directly aim for iteratively solving for v_* ? (\rightarrow value iteration)

Value iteration

1. Initialize $V(s) \in \mathbb{R}$ arbitrarily $\forall s \in \mathcal{S}$ except set $V(\text{terminal}) = 0$; set tolerance $\theta (=1e-6)$; set Δ to positive value
2. Find optimal value function
 - Loop while $\Delta > \theta$:
 - $\Delta \leftarrow 0$
 - Loop for each $s \in \mathcal{S}$
 - $v \leftarrow V(s)$
 - $V(s) = \max_{a \in \mathcal{A}} \sum_{s', r} [p(s', r | s, a) \cdot (R(s, a, s') + \gamma \cdot V(s'))]$
 - $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
3. Find corresponding optimal policy
 - $\pi(s) = \operatorname{argmax}_a \sum_{s', r} [p(s', r | s, a) \cdot (R(s, a, s') + \gamma \cdot V(s'))]$

UMassAmherst
The Commonwealth's Flagship Campus