Madhavan Mukund

https://www.cmi.ac.in/~madhavan

Advanced Machine Learning September–December 2021

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Set of states *S*, actions *A*, rewards *R*

э

▶ ∢ ⊒

→ ∢ Ξ

- Set of states S, actions A, rewards R
- At time t, agent in state S_t selects action A_t, moves to state S_{t+1} and receives reward R_{t+1}
 Trajectory S₀, A₀, R₁, S₁, A₁, R₂, S₂,...





- Set of states S, actions A, rewards R
- At time t, agent in state S_t selects action A_t, moves to state S_{t+1} and receives reward R_{t+1}
 Trajectory S₀, A₀, R₁, S₁, A₁, R₂, S₂,...
- Probabilistic transition function:
 p(s', r | s, a)
 - Probability of moving to state s' with reward r if we choose a at s
 - For each (s, a), $\sum_{s'} \sum_{r} p(s', r \mid s, a) = 1$



- Set of states S, actions A, rewards R
- At time t, agent in state S_t selects action A_t, moves to state S_{t+1} and receives reward R_{t+1}
 Trajectory S₀, A₀, R₁, S₁, A₁, R₂, S₂,...
- Probabilistic transition function: $p(s', r \mid s, a)$
 - Probability of moving to state s' with reward r if we choose a at s
 - For each (s, a), $\sum_{s'} \sum_{r} p(s', r \mid s, a) = 1$
 - Backup diagram



- Set of states *S*, actions *A*, rewards *R*
- At time t, agent in state S_t selects action A_t, moves to state S_{t+1} and receives reward R_{t+1}
 Trajectory S₀, A₀, R₁, S₁, A₁, R₂, S₂,...
- Probabilistic transition function: $p(s', r \mid s, a)$
 - Probability of moving to state s' with reward r if we choose a at s
 - For each (s, a), $\sum_{s'} \sum_{r} p(s', r \mid s, a) = 1$
 - Backup diagram
- Typically assume finite MDPs S, A and R are finite



Madhavan Mukund

MDP Example: Robot that collects empty cans

- State battery charge: high, low
- Actions: search for a can, wait for someone to bring can, recharge battery
 - No recharge when high



MDP Example: Robot that collects empty cans

- State battery charge: high, low
- Actions: search for a can, wait for someone to bring can, recharge battery
 - No recharge when high
- α, β, probabilities associated with change of battery state while searching



MDP Example: Robot that collects empty cans

- State battery charge: high, low
- Actions: search for a can, wait for someone to bring can, recharge battery
 - No recharge when high
- α, β, probabilities associated with change of battery state while searching
- 1 unit of reward per can collected
- r_{search} > r_{wait} cans collected while searching, waiting
- Negative reward for requiring rescue (low to high while searching)



How do we formalize long term rewards?

э

▶ ∢ ⊒

- How do we formalize long term rewards?
- Assume that each trajectory is a finite episode

э

- How do we formalize long term rewards?
- Assume that each trajectory is a finite episode
- Episode with T steps, expected reward at time t: $G_t \stackrel{\triangle}{=} R_{t+1} + R_{t+2} + \cdots + R_T$
 - Each episode is independent: rewards are reset after each episode

- How do we formalize long term rewards?
- Assume that each trajectory is a finite episode
- Episode with T steps, expected reward at time t: $G_t \stackrel{\triangle}{=} R_{t+1} + R_{t+2} + \cdots + R_T$
 - Each episode is independent: rewards are reset after each episode
- In some situations, trajectories may be (potentially) infinite

Discounted rewards: $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}, 0 \le \gamma \le 1$

Inductive calculation of expected reward

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

- How do we formalize long term rewards?
- Assume that each trajectory is a finite episode
- Episode with T steps, expected reward at time t: $G_t \stackrel{\triangle}{=} R_{t+1} + R_{t+2} + \cdots + R_T$
 - Each episode is independent: rewards are reset after each episode
- In some situations, trajectories may be (potentially) infinite

Discounted rewards: $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}, 0 \le \gamma \le 1$

Inductive calculation of expected reward

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

= $R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \gamma^{2} R_{t+3} + \cdots)$

- How do we formalize long term rewards?
- Assume that each trajectory is a finite episode
- Episode with T steps, expected reward at time t: $G_t \stackrel{\triangle}{=} R_{t+1} + R_{t+2} + \cdots + R_T$
 - Each episode is independent: rewards are reset after each episode
- In some situations, trajectories may be (potentially) infinite

Discounted rewards: $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}, 0 \le \gamma \le 1$

Inductive calculation of expected reward

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

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

Can make all episodes infinite by adding a self-loop with reward 0



э

-

Can make all episodes infinite by adding a self-loop with reward 0



• Allow $\gamma = 1$ only if sum converges

-

э

Can make all episodes infinite by adding a self-loop with reward 0



- Allow $\gamma = 1$ only if sum converges
- Alternatively, $G_t \stackrel{\triangle}{=} \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k$,

where we allow $T = \infty$ and $\gamma = 1$, but not both at the same time

Can make all episodes infinite by adding a self-loop with reward 0



- Allow $\gamma = 1$ only if sum converges
- Alternatively, $G_t \stackrel{\triangle}{=} \sum_{k=t+1}^{l} \gamma^{k-t-1} R_k$,

where we allow $T = \infty$ and $\gamma = 1$, but not both at the same time

If
$$T = \infty$$
, $R_k = +1$ for each k , $\gamma < 1$, then $G_t = \frac{1}{1 - \gamma}$

• A policy π describes how the agent chooses actions at a state

• $\pi(a \mid s)$ — probability of choosing *a* in state *s*, $\sum \pi(a \mid s) = 1$

• A policy π describes how the agent chooses actions at a state

•
$$\pi(a \mid s)$$
 — probability of choosing *a* in state *s*, $\sum_{a} \pi(a \mid s) = 1$

State value function at s, following policy π



• A policy π describes how the agent chooses actions at a state

•
$$\pi(a \mid s)$$
 — probability of choosing *a* in state *s*, $\sum_{a} \pi(a \mid s) = 1$

State value function at s, following policy π

$$v_{\pi}(s) \stackrel{ riangle}{=} \mathbb{E}_{\pi}[G_t \mid S_t = s] = \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s
ight]$$

• Action value function on choosing *a* at *s* and then following policy π

$$q_{\pi}(s,a) \stackrel{\Delta}{=} \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

• Note that $v_{\pi}(s) = \sum_{a} \pi(a \mid s) q_{\pi}(s, a)$

• A policy π describes how the agent chooses actions at a state

•
$$\pi(a \mid s)$$
 — probability of choosing *a* in state *s*, $\sum_{a} \pi(a \mid s) = 1$

State value function at s, following policy π

$$v_{\pi}(s) \stackrel{ riangle}{=} \mathbb{E}_{\pi}[G_t \mid S_t = s] = \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s
ight]$$

Action value function on choosing a at s and then following policy π

$$q_{\pi}(s,a) \stackrel{\triangle}{=} \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

• Note that
$$v_{\pi}(s) = \sum \pi(a \mid s)q_{\pi}(s, a)$$

а

Goal is to find an optimal policy, that maximizes state/action value at every state

•
$$v_{\pi}(s) \stackrel{\triangle}{=} \mathbb{E}_{\pi}[G_t \mid S_t = s]$$

3

・ロト ・ 同ト ・ ヨト ・ ヨ

•
$$v_{\pi}(s) \stackrel{\triangle}{=} \mathbb{E}_{\pi}[G_t \mid S_t = s]$$

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

3

イロト 不得 トイヨト イヨト

$$v_{\pi}(s) \stackrel{\triangle}{=} \mathbb{E}_{\pi}[G_{t} \mid S_{t} = s]$$

$$= \mathbb{E}_{\pi}(R_{t+1}) \cdot (\forall G_{t+1}) S_{t} = s]$$

$$= \sum_{a} \pi(a \mid s) \sum_{s'} \sum_{r} p(s', r \mid s, a)(r) \gamma \mathbb{E}_{\pi}[G_{t+1} \mid S_{t+1} = s']]$$

AML Sep-Dec 2021 7 / 16

3

・ロト ・ 同ト ・ ヨト ・ ヨ

•
$$v_{\pi}(s) \stackrel{\triangle}{=} \mathbb{E}_{\pi}[G_t \mid S_t = s]$$

 $= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} \mid S_t = s]$
 $= \sum_a \pi(a \mid s) \sum_{s'} \sum_r p(s', r \mid s, a) [r + \gamma \mathbb{E}_{\pi}[G_{t+1} \mid S_{t+1} = s']]$
 $= \sum_a \pi(a \mid s) \sum_{s'} \sum_r p(s', r \mid s, a) [r + \gamma v_{\pi}(s')]$

3

・ロト ・ 同ト ・ ヨト ・ ヨ

$$v_{\pi}(s) \stackrel{\Delta}{=} \mathbb{E}_{\pi}[G_{t} \mid S_{t} = s]$$

$$= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} \mid S_{t} = s]$$

$$= \sum_{a} \pi(a \mid s) \sum_{s'} \sum_{r} p(s', r \mid s, a) \left[r + \gamma \mathbb{E}_{\pi}[G_{t+1} \mid S_{t+1} = s'] \right]$$

$$= \sum_{a} \pi(a \mid s) \sum_{s'} \sum_{r} p(s', r \mid s, a) \left[r + \gamma v_{\pi}(s') \right]$$

Bellman equation relates state value at *s* to state values at successors of *s*

• Value function v_{π} is unique solution to the equation

Actions in each cell are {N,S,E,W}, with usual interpretation





э

▶ < ∃ ▶</p>

- Actions in each cell are $\{N,S,E,W\}$, with usual interpretation
- Reward is 0, except at boundaries
- Colliding with boundary position unchanged, reward -1



э

- Actions in each cell are {N,S,E,W}, with usual interpretation
- Reward is 0, except at boundaries
- Colliding with boundary position unchanged, reward -1
- Special squares A and B all four actions move as indicated, with rewards +10 and +5, respectively





- Actions in each cell are $\{N,S,E,W\}$, with usual interpretation
- Reward is 0, except at boundaries
- Colliding with boundary position unchanged, reward -1
- Special squares A and B all four actions move as indicated, with rewards +10 and +5, respectively
- Policy π choose each action with uniform probability 0.25





Madhavan Mukund

- Actions in each cell are {N,S,E,W}, with usual interpretation
- Reward is 0, except at boundaries
- Colliding with boundary position unchanged, reward -1
- Special squares A and B all four actions move as indicated, with rewards +10 and +5, respectively
- Policy π choose each action with uniform probability 0.25
- Solving Bellman equations, we obtain v_{π} for each square



0.1 0.7 0.7 0.4 -0.4 -1.0 -0.4 -0.4 -0.6 -1.2 -1.9 -1.3 -1.2 -1.4 -2.0

- Actions in each cell are {N,S,E,W}, with usual interpretation
- Reward is 0, except at boundaries
- Colliding with boundary position unchanged, reward -1
- Special squares A and B all four actions move as indicated, with rewards +10 and +5, respectively
- Policy π choose each action with uniform probability 0.25
- Solving Bellman equations, we obtain v_{π} for each square
- Values at boundary are negative



- Actions in each cell are {N,S,E,W}, with usual interpretation
- Reward is 0, except at boundaries
- Colliding with boundary position unchanged, reward -1
- Special squares A and B all four actions move as indicated, with rewards +10 and +5, respectively
- Policy π choose each action with uniform probability 0.25
- Solving Bellman equations, we obtain v_{π} for each square
- Values at boundary are negative
- Value at A is less than 10 because next move takes agent to boundary square with negative value







AML Sep-Dec 2021 8 / 16

- Actions in each cell are {N,S,E,W}, with usual interpretation
- Reward is 0, except at boundaries
- Colliding with boundary position unchanged, reward -1
- Special squares A and B all four actions move as indicated, with rewards +10 and +5, respectively
- Policy π choose each action with uniform probability 0.25
- Solving Bellman equations, we obtain v_{π} for each square
- Values at boundary are negative
- Value at A is less than 10 because next move takes agent to boundary square with negative value
- Value at B is more than 5 because next move is to a square with positive value





3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0
• Compare policies π , π' : $\pi \ge \pi'$ if $v_{\pi}(s) \ge v_{\pi'}(s)$ for every state s

- Compare policies π , π' : $\pi \geq \pi'$ if $v_{\pi}(s) \geq v_{\pi'}(s)$ for every state s
- Optimal policy π_* , $\pi_* \geq \pi$ for every π
 - Always exists, but may not be unique

- Compare policies π , π' : $\pi > \pi'$ if $v_{\pi}(s) > v_{\pi'}(s)$ for every state s
- Optimal policy $\pi_*, \pi_* \geq \pi$ for every π
 - Always exists, but may not be unique

• Optimal state value function, $v_*(s) \stackrel{\triangle}{=} \max_{\pi} v_{\pi}(s) = v_{\pi_*}(s)$

- Compare policies π , π' : $\pi \geq \pi'$ if $v_{\pi}(s) \geq v_{\pi'}(s)$ for every state s
- Optimal policy π_* , $\pi_* \geq \pi$ for every π
 - Always exists, but may not be unique
- Optimal state value function, $v_*(s) \stackrel{\triangle}{=} \max_{\pi} v_{\pi}(s) = v_{\pi_*}(s)$
- Optimal action value function, $q_*(s, a) \stackrel{\triangle}{=} \max_{\pi} q_{\pi}(s, a) = q_{\pi_*}(s, a)$

- Compare policies π , π' : $\pi \geq \pi'$ if $v_{\pi}(s) \geq v_{\pi'}(s)$ for every state s
- Optimal policy π_* , $\pi_* \geq \pi$ for every π
 - Always exists, but may not be unique
- Optimal state value function, $v_*(s) \stackrel{\triangle}{=} \max_{\pi} v_{\pi}(s) = v_{\pi_*}(s)$
- Optimal action value function, $q_*(s, a) \stackrel{\triangle}{=} \max_{\pi} q_{\pi}(s, a) = q_{\pi_*}(s, a)$
- Bellman optimality equation for v_*

 $v_*(s) = \max_a q_{\pi_*}(s,a)$

- Compare policies π , π' : $\pi \geq \pi'$ if $v_{\pi}(s) \geq v_{\pi'}(s)$ for every state s
- Optimal policy π_* , $\pi_* \geq \pi$ for every π
 - Always exists, but may not be unique
- Optimal state value function, $v_*(s) \stackrel{\triangle}{=} \max_{\pi} v_{\pi}(s) = v_{\pi_*}(s)$
- Optimal action value function, $q_*(s, a) \stackrel{\triangle}{=} \max_{\pi} q_{\pi}(s, a) = q_{\pi_*}(s, a)$
- Bellman optimality equation for v_*

$$egin{aligned} v_*(s) &= \max_a q_{\pi_*}(s,a) \ &= \max_a \mathbb{E}_{\pi_*}[G_t \mid S_t = s, A_t = a] \end{aligned}$$

- Compare policies π , π' : $\pi \geq \pi'$ if $v_{\pi}(s) \geq v_{\pi'}(s)$ for every state s
- Optimal policy π_* , $\pi_* \geq \pi$ for every π
 - Always exists, but may not be unique
- Optimal state value function, $v_*(s) \stackrel{\triangle}{=} \max_{\pi} v_{\pi}(s) = v_{\pi_*}(s)$
- Optimal action value function, $q_*(s, a) \stackrel{\triangle}{=} \max_{\pi} q_{\pi}(s, a) = q_{\pi_*}(s, a)$
- Bellman optimality equation for v_*

$$v_*(s) = \max_{a} q_{\pi_*}(s, a) \\ = \max_{a} \mathbb{E}_{\pi_*}[G_t \mid S_t = s, A_t = a] \\ = \max_{a} \mathbb{E}_{\pi_*}[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a]$$

- Compare policies π , π' : $\pi \geq \pi'$ if $v_{\pi}(s) \geq v_{\pi'}(s)$ for every state s
- Optimal policy π_* , $\pi_* \geq \pi$ for every π
 - Always exists, but may not be unique
- Optimal state value function, $v_*(s) \stackrel{\triangle}{=} \max_{\pi} v_{\pi}(s) = v_{\pi_*}(s)$
- Optimal action value function, $q_*(s,a) \stackrel{\triangle}{=} \max_{\pi} q_{\pi}(s,a) = q_{\pi_*}(s,a)$
- Bellman optimality equation for v_*

$$v_*(s) = \max_{a} q_{\pi_*}(s, a)$$

= $\max_{a} \mathbb{E}_{\pi_*}[G_t \mid S_t = s, A_t = a]$
= $\max_{a} \mathbb{E}_{\pi_*}[R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a]$
= $\max_{a} \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$

- Compare policies π , π' : $\pi \geq \pi'$ if $v_{\pi}(s) \geq v_{\pi'}(s)$ for every state s
- Optimal policy π_* , $\pi_* \geq \pi$ for every π
 - Always exists, but may not be unique
- Optimal state value function, $v_*(s) \stackrel{\triangle}{=} \max_{\pi} v_{\pi}(s) = v_{\pi_*}(s)$
- Optimal action value function, $q_*(s,a) \stackrel{\triangle}{=} \max_{\pi} q_{\pi}(s,a) = q_{\pi_*}(s,a)$
- Bellman optimality equation for v_*

$$v_{*}(s) = \max_{a} q_{\pi_{*}}(s, a)$$

= $\max_{a} \mathbb{E}_{\pi_{*}}[G_{t} | S_{t} = s, A_{t} = a]$
= $\max_{a} \mathbb{E}_{\pi_{*}}[R_{t+1} + \gamma G_{t+1} | S_{t} = s, A_{t} = a]$
= $\max_{a} \mathbb{E}[R_{t+1} + \gamma v_{*}(S_{t+1}) | S_{t} = s, A_{t} = a]$
= $\max_{a} \sum_{s', r} p(s', r | s, a)[r + \gamma v_{*}(s')]$

Madhavan Mukund

•
$$v_*(s) = \max_{a} \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

= $\max_{a} \sum_{s', r} p(s', r \mid s, a)[r + \gamma v_*(s')]$

•
$$v_*(s) = \max_{a} \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

= $\max_{a} \sum_{s', r} p(s', r \mid s, a)[r + \gamma v_*(s')]$

Likewise, for action value function

$$q_*(s, a) = \mathbb{E}[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = t, A_t = a]$$

= $\sum_{s', r} p(s', r \mid s, a)[r + \max_{a'} \gamma q_*(s', a')]$

•
$$v_*(s) = \max_{a} \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

= $\max_{a} \sum_{s', r} p(s', r \mid s, a)[r + \gamma v_*(s')]$

Likewise, for action value function

$$q_*(s, a) = \mathbb{E}[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = t, A_t = a]$$

= $\sum_{s', r} p(s', r \mid s, a)[r + \max_{a'} \gamma q_*(s', a')]$

- For finite state MDPs, can solve explicitly for v_*
 - **n** states, *n* equations in *n* unknowns, (assuming we know p)

•
$$v_*(s) = \max_{a} \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

= $\max_{a} \sum_{s', r} p(s', r \mid s, a)[r + \gamma v_*(s')]$

Likewise, for action value function

$$q_*(s, a) = \mathbb{E}[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = t, A_t = a]$$

= $\sum_{s', r} p(s', r \mid s, a)[r + \max_{a'} \gamma q_*(s', a')]$

- For finite state MDPs, can solve explicitly for v_*
 - **n** states, *n* equations in *n* unknowns, (assuming we know p)
- However, *n* is usually large, computationally infeasible
 - State space of a game like chess or Go

•
$$v_*(s) \neq \max_{a} \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a]$$

= $\max_{a} \sum_{s',r} p(s',r | s,a)[r + \gamma v_*(s')]$

Likewise, for action value function

$$q_*(s, a) = \mathbb{E}[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = t, A_t = a]$$

= $\sum_{s', r} p(s', r \mid s, a)[r + \max_{a'} \gamma q_*(s', a')]$

- For finite state MDPs, can solve explicitly for v_*
 - **n** states, *n* equations in *n* unknowns, (assuming we know p)
- However, n is usually large, computationally infeasible
 - State space of a game like chess or Go
- \blacksquare Instead, we will explore iterative methods to approximate v_{\ast}

Madhavan Mukund

• Given a policy π , compute its state value function v_{π}

- Given a policy π , compute its state value function v_{π}
- Bellman equations: $v_{\pi}(s) = \sum_{a} \pi(a \mid s) \sum_{s'} \sum_{r} p(s', r \mid s, a) [r + \gamma v_{\pi}(s')]$
 - For MDP with n states, n equations in n unknowns
 - Can solve to get v_{π} , but computationally infeasible for large *n*

- Given a policy π , compute its state value function v_{π}
- Bellman equations: $v_{\pi}(s) = \sum_{a} \pi(a \mid s) \sum_{s'} \sum_{r} p(s', r \mid s, a) [r + \gamma v_{\pi}(s')]$
 - For MDP with n states, n equations in n unknowns
 - Can solve to get v_{π} , but computationally infeasible for large n
- Instead, use the Bellman equations as update rules

- Given a policy π , compute its state value function v_{π}
- Bellman equations: $v_{\pi}(s) = \sum_{a} \pi(a \mid s) \sum_{s'} \sum_{r} p(s', r \mid s, a) [r + \gamma v_{\pi}(s')]$
 - For MDP with *n* states, *n* equations in *n* unknowns
 - Can solve to get v_{π} , but computationally infeasible for large n
- Instead, use the Bellman equations as update rules
 - Initialize $v_{\pi}^{0}(s)$: set $v_{\pi}^{0}(\text{term}) = 0$ for terminal state term, arbitrary values for other s

- Given a policy π , compute its state value function v_{π}
- Bellman equations: $v_{\pi}(s) = \sum_{a} \pi(a \mid s) \sum_{s'} \sum_{r} p(s', r \mid s, a) [r + \gamma v_{\pi}(s')]$
 - For MDP with n states, n equations in n unknowns
 - Can solve to get v_{π} , but computationally infeasible for large *n*
- Instead, use the Bellman equations as update rules
 - Initialize $v_{\pi}^{0}(s)$: set $v_{\pi}^{0}(\text{term}) = 0$ for terminal state term, arbitrary values for other s
 - Update v_{π}^k to v_{π}^{k+1} using: $v_{\pi}^{k+1}(s) = \sum_{a} \pi(a \mid s) \sum_{s'} \sum_{r} p(s', r \mid s, a) \left[r + \gamma v_{\pi}^k(s') \right]$

- Given a policy π , compute its state value function v_{π}
- Bellman equations: $v_{\pi}(s) = \sum_{a} \pi(a \mid s) \sum_{s'} \sum_{r} p(s', r \mid s, a) [r + \gamma v_{\pi}(s')]$
 - For MDP with n states, n equations in n unknowns
 - Can solve to get v_{π} , but computationally infeasible for large *n*
- Instead, use the Bellman equations as update rules
 - Initialize $v_{\pi}^{0}(s)$: set $v_{\pi}^{0}(\text{term}) = 0$ for terminal state term, arbitrary values for other s
 - Update v_{π}^k to v_{π}^{k+1} using: $v_{\pi}^{k+1}(s) = \sum_{a} \pi(a \mid s) \sum_{s'} \sum_{r} p(s', r \mid s, a) [r + \gamma v_{\pi}^k(s')]$
 - Stop when incremental change $\Delta = |v_{\pi}^{k+1} v_{\pi}^{k}|$ is below threshold θ

- Assume a deterministic policy π
- Using v_{π} , can we find a better policy π' ?

- Assume a deterministic policy π
- Using v_{π} , can we find a better policy π' ?
- Is there a state s where we can substitute $\pi(s)$ by a better choice a?

- Assume a deterministic policy π
- Using v_{π} , can we find a better policy π' ?
- Is there a state s where we can substitute $\pi(s)$ by a better choice a?

$$q_{\pi}(s,a) = \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = a]$$
$$= \sum_{s',r} p(s',r \mid s,a) [r + \gamma v_{\pi}(s')]$$

- Assume a deterministic policy π
- Using v_{π} , can we find a better policy π' ?
- Is there a state s where we can substitute $\pi(s)$ by a better choice a?

•
$$q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = a]$$

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

• If $q_{\pi}(s, a) > v_{\pi}(s)$, modify π so that $\pi(s) = a$

- Assume a deterministic policy π
- Using v_{π} , can we find a better policy π' ?
- Is there a state s where we can substitute $\pi(s)$ by a better choice a?

•
$$q_{\pi}(s, a) = \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = a]$$

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

- If $q_{\pi}(s, a) > v_{\pi}(s)$, modify π so that $\pi(s) = a$
- The new policy π' is strictly better

For deterministic policies π , π' :

- If $q_{\pi}(s,\pi'(s)) \geq v_{\pi}(s)$ for all s, then $\pi' \geq \pi$,
- If $\pi' \ge \pi$ and $q_{\pi}(s, \pi'(s)) > v_{\pi}(s)$ for some s, then $v_{\pi'}(s) > v_{\pi}(s)$.

3

1 E N

For deterministic policies π , π' :

- If $q_{\pi}(s,\pi'(s)) \geq v_{\pi}(s)$ for all s, then $\pi' \geq \pi$,
- If $\pi' \ge \pi$ and $q_{\pi}(s, \pi'(s)) > v_{\pi}(s)$ for some s, then $v_{\pi'}(s) > v_{\pi}(s)$.

Proof of the theorem is not difficult for deterministic policies

For deterministic policies π , π' :

- If $q_{\pi}(s,\pi'(s)) \geq v_{\pi}(s)$ for all s, then $\pi' \geq \pi$,
- If $\pi' \ge \pi$ and $q_{\pi}(s, \pi'(s)) > v_{\pi}(s)$ for some s, then $v_{\pi'}(s) > v_{\pi}(s)$.
- Proof of the theorem is not difficult for deterministic policies
- The theorem extends to probabilistic policies also

For deterministic policies π , π' :

- If $q_{\pi}(s,\pi'(s)) \geq v_{\pi}(s)$ for all s, then $\pi' \geq \pi$,
- If $\pi' \ge \pi$ and $q_{\pi}(s, \pi'(s)) > v_{\pi}(s)$ for some s, then $v_{\pi'}(s) > v_{\pi}(s)$.

- Proof of the theorem is not difficult for deterministic policies
- The theorem extends to probabilistic policies also
- Provides a basis to iteratively improve the policy

Start with a random policy π_0

3

▶ ∢ ⊒

< □ > < 同

- Start with a random policy π_0
- Use policy evaluation to compute v_{π_0}

- Start with a random policy π_0
- Use policy evaluation to compute v_{π_0}
- Use policy improvement to construct a better policy π_1

- Start with a random policy π_0
- Use policy evaluation to compute v_{π_0}
- Use policy improvement to construct a better policy π_1
- Policy iteration: Alternate between policy evaluation and policy improvement $\pi_0 \xrightarrow{\text{evaluate}} v_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} v_{\pi_1} \xrightarrow{\text{improve}} \pi_2 \xrightarrow{\text{evaluate}} \cdots$

- Start with a random policy π_0
- Use policy evaluation to compute v_{π_0}
- Use policy improvement to construct a better policy π_1
- Policy iteration: Alternate between policy evaluation and policy improvement $\pi_0 \xrightarrow{\text{evaluate}} v_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} v_{\pi_1} \xrightarrow{\text{improve}} \pi_2 \xrightarrow{\text{evaluate}} \cdots \xrightarrow{\text{improve}} \pi_* \xrightarrow{\text{evaluate}} v_{\pi_*}$
- Finite MDPs can improve π only finitely many times,
 - Must converge to optimal policy

- Start with a random policy π_0
- Use policy evaluation to compute v_{π_0}
- Use policy improvement to construct a better policy π_1
- Policy iteration: Alternate between policy evaluation and policy improvement $\pi_0 \xrightarrow{\text{evaluate}} v_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} v_{\pi_1} \xrightarrow{\text{improve}} \pi_2 \xrightarrow{\text{evaluate}} \cdots \xrightarrow{\text{improve}} \pi_* \xrightarrow{\text{evaluate}} v_{\pi_*}$
- Finite MDPs can improve π only finitely many times,
 - Must converge to optimal policy
- Nested iteration each policy evaluation is itself an iteration
 - Speed up by using v_{π_i} as initial state to compute $v_{\pi_{i+1}}$

= nar

Value iteration

Policy iteration — policy evaluation requires a nested iteration
- Policy iteration policy evaluation requires a nested iteration
- A partial computation of v_{π_k} is sufficent to proceed towards π_* , v_*

э

- Policy iteration policy evaluation requires a nested iteration
- A partial computation of v_{π_k} is sufficent to proceed towards π_* , v_*
- Even a single iteration in the computation of v_{π_k} will do

э

- Policy iteration policy evaluation requires a nested iteration
- A partial computation of v_{π_k} is sufficent to proceed towards π_* , v_*
- Even a single iteration in the computation of v_{π_k} will do
- Combine policy improvement and one step update at each state

- Policy iteration policy evaluation requires a nested iteration
- A partial computation of v_{π_k} is sufficent to proceed towards π_* , v_*
- Even a single iteration in the computation of v_{π_k} will do
- Combine policy improvement and one step update at each state
- Value iteration

$$\begin{aligned} v_{\pi_{k+1}}(s,a) &= \max_{a} \mathbb{E}[R_{t+1} + \gamma v_{\pi_k}(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_{a} \sum_{s',r} p(s',r \mid s,a) \left[r + \gamma v_{\pi_k}(s')\right] \end{aligned}$$

- Policy iteration policy evaluation requires a nested iteration
- A partial computation of v_{π_k} is sufficent to proceed towards π_* , v_*
- Even a single iteration in the computation of v_{π_k} will do
- Combine policy improvement and one step update at each state
- Value iteration

$$egin{aligned} & \mathcal{V}_{\pi_{k+1}}(s,a) = \max_{a} \mathbb{E}[R_{t+1} + \gamma \mathbf{v}_{\pi_k}(S_{t+1}) \mid S_t = s, A_t = a] \ & = \max_{a} \sum_{s',r} p(s',r \mid s,a) \left[r + \gamma \mathbf{v}_{\pi_k}(s')
ight] \end{aligned}$$

• Again, stop when incremental change $\Delta = |v_{\pi_{k+1}} - v_{\pi_k}|$ is below threshold θ

In the literature, policy iteration and value iteration are referred to as dynamic programming methods

3

▶ < ∃ ▶</p>

< □ > < 向

- In the literature, policy iteration and value iteration are referred to as dynamic programming methods
- Requires knowledge of the model $-p(s', r \mid s, a)$

3

4 E N

- In the literature, policy iteration and value iteration are referred to as dynamic programming methods
- Requires knowledge of the model $-p(s', r \mid s, a)$
- How to combine policy evaluation and policy improvement is flexible
 - Value iteration is policy iteration with policy evaluation truncated to a single step
 - Generalized policy iteration simultaneously maintain and update approximations of π_* and v_*