Policy and Value Iteration

Madhavan Mukund

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

Advanced Machine Learning September–December 2021

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで



AML Sep–Dec 2021 2 / 11

- 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*

Given II -> compile

Instead, use the Bellman equations as update rules

(*+4)-(4,*)

 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
- 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: $v_{\pi}^{k+1}(s) = \sum \pi(a \mid s) \sum \sum p(s', r \mid s, a) \left[r + (v_{\pi}^{k})^{s'} \right]$



- 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: $v_{\pi}^{k+1}(s) = \sum_{a} \pi(a \mid s) \sum_{s'} \sum_{r} p(s', r \mid s, a) [r + \gamma v_{\pi}^k(s')]$

- 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

Iterative Policy Evaluation, for estimating $V \approx v_{\pi}$

```
Input \pi, the policy to be evaluated
Algorithm parameter: a small threshold \theta > 0 determining accuracy of estimation
Initialize V(s), for all s \in S^+, arbitrarily except that V(terminal) = 0
  Dop:

△ ← 0 - max difference (absorbule) deeved in this
Loop:
  Loop for each s \in S:
       v \leftarrow V(s)
       V(s) \leftarrow \sum_{a} \pi(a|s) \sum_{s' r} p(s', r|s, a) [r + \gamma V(s')]
      \Delta \leftarrow \max(\Delta, |v - V(s)|)
until \Delta < \theta
```

▶ < ⊒ ▶

Policy evaluation example







- 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?

TL(a(s) - 1 for closen a 1) othermore P(s',r(s,a) - prolablade

- 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

T1 = TT except at s

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 $v_{\pi'}(s) > v_{\pi}(s)$.

э

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

э

▶ ∢ ⊒

< □ > < □ > < □ > < Ξ

- 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 $\pi_0 \xrightarrow{\text{evaluate}} v_{\pi_0} \xrightarrow{\text{improve}} \pi_1 \xrightarrow{\text{evaluate}} v_{\pi_1} \xrightarrow{\text{improve}} \pi_2 \xrightarrow{\text{evaluate}} \cdots$

4 E N

- 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 $v_{\pi_{i+1}}$



Optimizing Policy Iteration



0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0



0.0 -14. -20.

-14. -18. -20.

-20. -20. -18.

-22. -20. -14. 0.0

-22.

-20

-14.



k = 3

э

イロト 不得 トイヨト イヨト

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_*

- 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

- 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

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

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

I step value update

- 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



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

э

2 E

- 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_*

- 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_*
- Asynchronous dynamic programming for large state spaces

Fair exploration of entire space

3