第四课 value iteration & policy iteration

前排提醒,感觉这节课没啥用,而且挺绕

1.value iteration, 也就是求解BOE的算法

Value iteration algorithm

The algorithm

$$v_{k+1} = f(v_k) = \max_{\pi} (r_{\pi} + \gamma P_{\pi} v_k), \quad k = 1, 2, 3...$$

can be decomposed to two steps.

• Step 1: policy update. This step is to solve

$$(\pi_{k+1}) = \arg\max_{\pi} (r_{\pi} + \gamma P_{\pi} v_k)$$

where v_k is given.

Step 2: value update.

$$v_{k+1} = r_{\overline{x_{k+1}}} + \gamma P_{\overline{x_{k+1}}} v_k$$

Question: is v_k a state value? No, because it is not ensured that v_k satisfies a Bellman equation.

Value iteration algorithm - Pseudocode

▷ Procedure summary:

Vo

$$v_k(s) o \underline{q_k(s,a)} o ext{greedy policy} \underline{\pi_{k+1}}(a|s) o ext{new value} \underline{v_{k+1}} = \max_a q_k(s,a)$$

Pseudocode: Value iteration algorithm

Initialization: The probability model p(r|s,a) and p(s'|s,a) for all (s,a) are known. Initial guess v_0 .

Aim: Search the optimal state value and an optimal policy solving the Bellman optimality equation.

While v_k has not converged in the sense that $||v_k - v_{k-1}||$ is greater than a predefined small threshold, for the kth iteration, do

For every state $s \in \mathcal{S}$, do

For every action $a \in \mathcal{A}(s)$, do

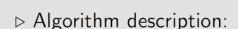
q-value:
$$q_k(s,a) = \sum_r p(r|s,a)r + \gamma \sum_{s'} p(s'|s,a)v_k(s')$$

Maximum action value: $a_k^*(s) = \arg\max_a q_k(a, s)$

Policy update: $\pi_{k+1}(a|s) = 1$ if $a = a_k^*$ and $\pi_{k+1}(a|s) = 0$ otherwise

 \longrightarrow Value update: $v_{k+1}(s) = \max_{a} q_k(a, s)$

Policy iteration algorithm



Given a random initial policy π_0 ,



This step is to calculate the state value of π_k

$$v_{\pi_k} = r_{\pi_k} + \gamma P_{\pi_k} v_{\pi_k} \longleftarrow$$

김 🕞

Note that v_{π_k} is a state value function.

Step 2: policy improvement (PI)

$$(\pi_{k+1}) = \arg\max_{\pi} (r_{\pi} + \gamma P_{\pi} (r_{\pi_k}))$$

The maximization is componentwise!

这里policy evaluation是通过求解BE的递归算法得到的

Policy iteration algorithm - Implementation

Pseudocode: Policy iteration algorithm

Initialization: The probability model p(r|s,a) and p(s'|s,a) for all (s,a) are known. Initial guess π_0 .

Aim: Search for the optimal state value and an optimal policy.

While the policy has not converged, for the kth iteration, do

Initialization: an arbitrary initial guess $v_{\pi_k}^{(0)}$

While $v_{\pi_k}^{(j)}$ has not converged, for the jth iteration, do

For every state
$$s \in \mathcal{S}$$
, do

$$v_{\pi_k}^{(j+1)}(s) = \sum_{a} \pi_k(a|s) \left[\sum_{r} p(r|s, a)r + \gamma \sum_{s'} p(s'|s, a)v_{\pi_k}^{(j)}(s') \right]$$

Policy improvement:

For every state $s \in \mathcal{S}$, do

For every action $a \in \mathcal{A}(s)$, do

$$q_{\pi_k}(s, a) = \sum_r p(r|s, a)r + \gamma \sum_{s'} p(s'|s, a)v_{\pi_k}(s')$$

$$a_k^*(s) = \arg\max_a q_{\pi_k}(s, a)$$

$$a_k^*(s) = \arg\max_a q_{\pi_k}(s,a)$$

$$\pi_{k+1}(a|s) = 1 \text{ if } a = a_k^*, \text{ and } \pi_{k+1}(a|s) = 0 \text{ otherwise}$$

Compare value iteration and policy iteration

▷ Let's compare the steps carefully:

| | Policy iteration algorithm | Value iteration algorithm | Comments |
|------------|---|--|--|
| 1) Policy: | π_0 | N/A | |
| 2) Value: | $v_{\pi_0} = r_{\pi_0} + \gamma P_{\pi_0} v_{\pi_0}$ | $v_0 := v_{\pi_0}$ | |
| 3) Policy: | $\pi_1 = \arg\max_{\pi} (r_{\pi} + \gamma P_{\pi} v_{\pi_0})$ | $\pi_1 = \arg\max_{\pi} (r_{\pi} + \gamma P_{\pi} v_0)$ | The two policies are the |
| | | | same |
| 4) Value: | $v_{\pi_1} = r_{\pi_1} + \gamma P_{\pi_1} v_{\pi_1}$ | $v_1 = r_{\pi_1} + \gamma P_{\pi_1} v_0$ | $v_{\pi_1} \geq v_1 \text{ since } v_{\pi_1} \geq$ |
| | | | v_{π_0} |
| 5) Policy: | $\pi_2 = \arg\max_{\pi} (r_{\pi} + \gamma P_{\pi} v_{\pi_1})$ | $\pi_2' = \arg\max_{\pi} (r_{\pi} + \gamma P_{\pi} v_1)$ | |
| : | : | : | : |
| | • | • | • |

- They start from the same initial condition.
- The first three steps are the same.
- The fourth step becomes different:
 - In policy iteration, solving $v_{\pi_1} = r_{\pi_1} + \gamma P_{\pi_1} v_{\pi_1}$ requires an iterative algorithm (an infinite number of iterations)
 - In value iteration, $v_1 = r_{\pi_1} + \gamma P_{\pi_1} v_0$ is a one-step iteration

Compare value iteration and policy iteration

Consider the step of solving $v_{\pi_1} = r_{\pi_1} + \gamma P_{\pi_1} v_{\pi_1}$:

$$\begin{aligned} v_{\pi_1}^{(0)} &= v_0 \\ \text{value iteration} \leftarrow v_1 &\longleftarrow v_{\pi_1}^{(1)} &= r_{\pi_1} + \gamma P_{\pi_1} v_{\pi_1}^{(0)} \\ v_{\pi_1}^{(2)} &= r_{\pi_1} + \gamma P_{\pi_1} v_{\pi_1}^{(1)} \\ &\vdots \\ \text{truncated policy iteration} \leftarrow \bar{v}_1 &\longleftarrow v_{\pi_1}^{(j)} &= r_{\pi_1} + \gamma P_{\pi_1} v_{\pi_1}^{(j-1)} \\ &\vdots \end{aligned}$$

policy iteration
$$\leftarrow v_{\pi_1} \leftarrow v_{\pi_1}^{(\infty)} = r_{\pi_1} + \gamma P_{\pi_1} v_{\pi_1}^{(\infty)}$$

- The value iteration algorithm computes once.
- The policy iteration algorithm computes an infinite number of iterations.
- The truncated policy iteration algorithm computes a finite number of shiyu Zhao iterations (say j). The rest iterations from j to ∞ are truncated.