

Contraction Mapping Theorem

Theorem (Contraction Mapping Theorem)  
For any equation that has the form of  $x = f(x)$ , if  $f$  is a contraction mapping, then

- Existence: there exists a fixed point  $x^*$  satisfying  $f(x^*) = x^*$ .
- Uniqueness: The fixed point  $x^*$  is unique.
- Algorithm: Consider a sequence  $\{x_k\}$  where  $x_{k+1} = f(x_k)$ , then  $x_k \rightarrow x^*$  as  $k \rightarrow \infty$ . Moreover, the convergence rate is exponentially fast.

value iteration algorithm

algorithm to solve the problem

existence

uniqueness

optimality

$v^* = r_{\pi^*} + \gamma P_{\pi^*} v^*$

所有reward进行线性变化ar+b, policy不改变

改变r: reward

改变γ: discount rate

γ小, 近视

definition

analyzing

Optimal policy

Optimal state value

找到

step1: policy update

step2: value update

value iteration

policy iteration

truncated policy iteration

3. BOE: Bellman Optimal Equation

$$v = \max_{\pi}(r_{\pi} + \gamma P_{\pi} v)$$

BOE

value iteration & policy iteration

RL: Find Optimal Policy

2. The Bellman equation

1. Markov Decision Process

action value

$q_{\pi}(s, a) = \mathbb{E}[G_t | S_t = s, A_t = a]$

policy evaluation

closed form实际不用, 逆计算量大

iterative solution 从V0开始, 迭代收敛

$$v_{\pi} = (I - \gamma P_{\pi})^{-1} r_{\pi}$$
$$v_{k+1} = r_{\pi} + \gamma P_{\pi} v_k$$

state value

$$v_{\pi}(s) = \mathbb{E}[G_t | S_t = s]$$

return

bootstrapping

Method 2:

s1

s2

s3

s4

$v_1$

$v_2$

$v_3$

$v_4$

$v_1 = r_1 + \gamma(r_2 + \gamma r_3 + \dots) = r_1 + \gamma v_2$   
 $v_2 = r_2 + \gamma(r_3 + \gamma r_4 + \dots) = r_2 + \gamma v_3$   
 $v_3 = r_3 + \gamma(r_4 + \gamma r_1 + \dots) = r_3 + \gamma v_4$   
 $v_4 = r_4 + \gamma(r_1 + \gamma r_2 + \dots) = r_4 + \gamma v_1$

①  
②

which can be rewritten as

$$\mathbf{v} = \mathbf{r} + \gamma \mathbf{P} \mathbf{v}$$

Deterministic situation

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

mean of immediate rewards      mean of future rewards

$$= \sum_a \pi(a|s) \left[ \sum_r p(r|s, a) r + \gamma \sum_{s'} p(s'|s, a) v_{\pi}(s') \right], \quad \forall s \in \mathcal{S}$$

Highlights: symbols in this equation

- $v_{\pi}(s)$  and  $v_{\pi}(s')$  are state values to be calculated. Bootstrapping!
- $\pi(a|s)$  is a given policy. Solving the equation is called policy evaluation.
- $p(r|s, a)$  and  $p(s'|s, a)$  represent the dynamic model. What if the model is known or unknown?

General form