

Value iteration

$$V(S) = \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s)[r(s, a, s') + \gamma V(s')]$$

Algorithm

Value iteration:

- 1) Set V_0 to arbitrary value for each s in S (choose 0 as the value)
- 2) While diff is $\geq \epsilon$ do
 - a. For each s in S do
 - i. $V_{t+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s)[r(s, a, s') + \gamma V_t(s')]$
- 3) Select policy

Value Iteration

$$V_{t+1}(s) :=$$

$$\max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s)[r(s, a, s') + \gamma V_t(s')]$$

0	$V_1 = 0$ $V_2 = 0.5$	$V_1 = 0.72$ $V_2 = 0.72$ 0.7248	+1
0	0.8	$V_1 = 0$ $V_2 = ?$	-1
0	0	0.1	0

Assuming $\gamma = 0.9$

$$v_2 = 0.8*(0 + 0.9*1) + 0.1*(0 + 0.9*0) + 0.1*(0 + 0.9*0) = 0.72$$

$$0.8*(0 + 0.9*0) + \dots$$

0 ✓

$$0.8*(0 + 0.9*0) = 0.72$$

Deciding how to act

$$0.8*(0 + 0.9*1) + 0.1*(0 + 0.9*0)$$

$$+ 0.1*(0 + 0.9*0.72) = 0.7848$$

$$\underset{a \in A}{\operatorname{Argmax}} \sum_{s' \in S} P_a(s'|s) [r(s, a, s') + \gamma V(s')]$$

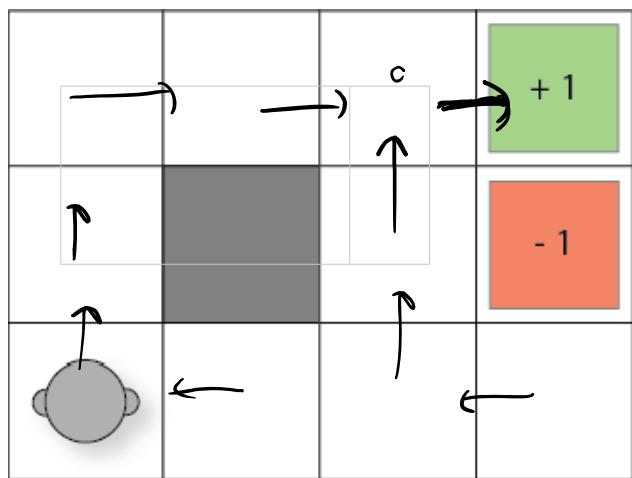
$Q(s, a)$

Just intuition in P0054

Policy iteration

$$0.6 * (2 + 0.9 * 10) = 6.6$$

$$0.4 * (0 + 0.9 * 12) = 4.32$$



π (random assign initially)

①

Policy evaluation (How good is this state here if taking the given action)

②

Policy update

choose the max (Best) action

VI

$$O(|S|^2 |A|)$$

PI

$$O(|S|^4 |A| + |A|^3)$$

$$V^\pi(s) = \sum_{s' \in S} P(s'|s) [r(s, a, s') + \gamma V^\pi(s')] \quad a = \pi(s)$$