

## Value iteration

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

Value iteration:

- 1) Set  $V_0$  to arbitrary value for each  $s$  in  $S$  (choose 0 as the value)
- 2) While  $\text{diff} \geq \epsilon$ 
  - 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

$$V_{t+1}(s) := \max_{a \in A(s)} \sum_{s' \in S} P_a(s'|s) [r(s, a, s') + \gamma V_t(s')]$$

○	○	$V_1 = 0$ $V_2 = 0.72$ $\uparrow 0.1$ $V_1 = 0$ $V_2 = ?$	+1
○	0.8	$\leftarrow 0.8$ $\downarrow 0.1$ $V_1 = ?$	-1
○	○	$\checkmark$ $0.1$ $\downarrow 0.1$	○

Assuming  $\gamma = 0.9$

$$V_2 = 0.8 * (0 + 0.9 * 1) + 0.1 * (0 + 0.9 * 0) + 0.1 * (0 + 0.9 * 0) = 0.72$$

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

○ ✓

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

$\bigcirc$   $V_1$   $\rightarrow$   $V_2$

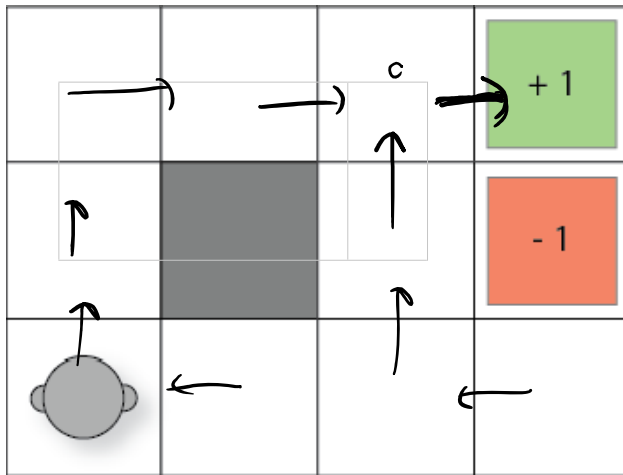
Deciding how to act

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

# Policy iteration

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

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



$\pi$

Policy evaluation  
Policy update

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

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

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

$a = \pi(s)$