# Reinforcement Learning Homework 01

Qiu Yihang

Feb 2023

## 1 Convergence of Policy Iteration

*Proof.* Let the state space be $\mathcal{S}$, the action space be $\mathcal{A}$.

To prove the new policy produced by policy iterations will be at least as good as the original one, we need to prove

$$V^{\pi_{i+1}}(s) \geq V^{\pi_i}(s).$$

For $V^{\pi_i}(s)$, by the process of policy iteration, we know

$$\sum_{s'} P(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V^{\pi_i}(s')] \leq \sum_{s'} \max_{a \in \mathcal{A}} P(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V^{\pi_i}(s')]$$

Thus,

$$
\begin{aligned}
V^{\pi_i}(s) \leq Q(s, &\pi_{i+1}(s)) \\
&= \mathbb{E}_{\pi_{i+1}} [R(s_0, \pi_{i+1}(s_0), s_1) + \gamma V^{\pi_i}(s_1) \mid \pi_i, s_0 = s] \\
&\leq \mathbb{E}_{\pi_{i+1}} [R(s_0, \pi_{i+1}(s_0), s_1) + \gamma Q(s_1, \pi_{i+1}(s_1)) \mid \pi_i, s_0 = s] \\
&= \mathbb{E}_{\pi_{i+1}} [R(s_0, \pi_{i+1}(s_0), s_1) + \gamma R(s_1, \pi_{i+1}(s_1), s_2) + \gamma^2 \cdot V^{\pi_i}(s_1) \mid \pi_i, s_0 = s] \\
&\leq \mathbb{E}_{\pi_{i+1}} [R(s_0, \pi_{i+1}(s_0), s_1) + \gamma R(s_1, \pi_{i+1}(s_1), s_2) + \gamma^2 Q(s_2, \pi_{i+1}(s_2)) \mid \pi_i, s_0 = s] \\
&= \dots \\
&\leq \mathbb{E} \left[ \sum_{t=0}^{\infty} R(s_t, \pi_{i+1}(s_t), s_{t+1}) \right] = V^{\pi_{i+1}}(s)
\end{aligned}
\tag{1}
$$

Moreover, $V^{\pi_{i+1}}(s) = V^{\pi_i}(s)$ **iff.**

$$
\forall s' \in \mathcal{S}, \ \max_a P(s, a, s') [R(s, a, s') + \gamma V^{\pi_i}(s')] = P(s, \pi_i(s), s') [R(s, \pi_i(s), s') + \gamma V^{\pi_i}(s')]
$$

$$
\text{i.e. } \pi_i(s) \in \arg\max_a \sum_{s'} Q^{\pi_i}(s, a), \ \text{i.e. } \pi_i(s) \text{ is already an optimal policy.}
\tag{2}
$$

Thus, the new policy is at least as good as the original one. $\square$

We prove the convergence of policy iterations as follows.

Let the set of optimal policies be $\Pi^*$.

Given that the state space, the action space and the reward function is all finite, we have

$$\forall s, s' \in \mathcal{S}, a \in \mathcal{A}, \ R_- \le R(s, a, s') < R_+.$$

Then for any possible policy $\pi$, since $0 \le \gamma < 1$,

$$\frac{R_-}{1-\gamma} \le V^\pi(s) < \sum_{t=0}^{\infty} \gamma^t \cdot R_0 = \frac{R_+}{1-\gamma}$$

By 1 and 2, we know

$$\begin{cases} V^{\pi_i} = V^{\pi_{i+1}} & \textbf{iff.} \ \forall s \in \mathcal{S}, \pi_i(s), \pi_{i+1}(s) \in \arg\max_a \sum_{s'} Q^{\pi_i}(s, a). \\ \exists s \in \mathcal{S}, V^{\pi_i}(s) < V^{\pi_{i+1}}(s) & \text{otherwise} \end{cases}$$

i.e.

$$\begin{cases} V^{\pi_i} = V^{\pi_{i+1}} & \textbf{iff.} \ \pi_i, \pi_{i+1} \in \Pi^*. \\ \exists s \in \mathcal{S}, V^{\pi_i}(s) < V^{\pi_{i+1}}(s) & \text{otherwise} \end{cases}$$

Meanwhile, there is at most $|\mathcal{S}| \times |\mathcal{A}|$ polices. Then for $N \ge |\mathcal{S}| \times |\mathcal{A}| + 1$, we know

$$\forall s \in \mathcal{S}, \exists k_s \in \mathbb{N}, V^{\pi_1}(s) < V^{\pi_2}(s) < ... < V^{\pi_{k_s-1}}(s) = V^{\pi_{k-s}}(s) = ... = V^{\pi_N}(s).$$

i.e.

$$V^{\pi_{N-1}} = V^{\pi_N} \implies \pi_{N-1}(s) \in \Pi^*, \pi_N(s) \in \Pi^*.$$

Moreover, by the definition of how $\pi_{i+1}$ derives from $\pi_i$, we know when $\pi_N \in \Pi^*$, $\pi_{N+1} = \pi_N$.

Therefore, when $i \to \infty, \pi_i \to \pi_N \in \Pi^*$.

Thus, the policy iteration converges to an optimal policy. ∎
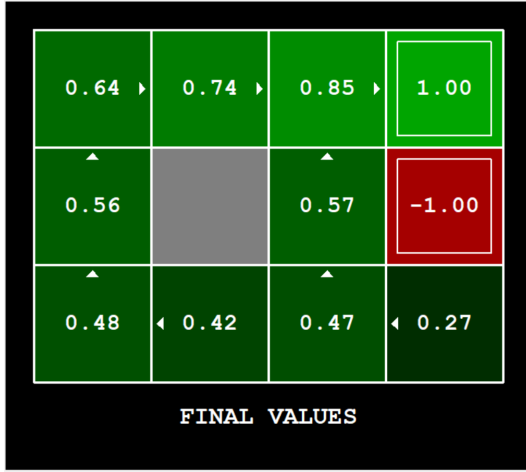
# 2 Grid World

## 2.1 Value Iteration

The value iteration is implemented in `valueIterationAgent.py`.

The results given by value iteration under $\varepsilon = 0.01$ or $0.001$ are shown in Fig. 1.
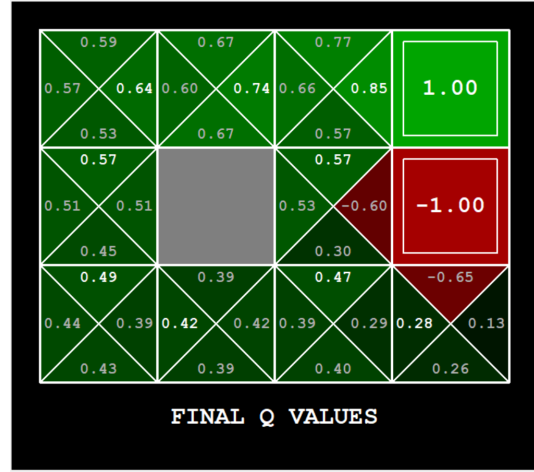
## 2.2 Policy Iteration

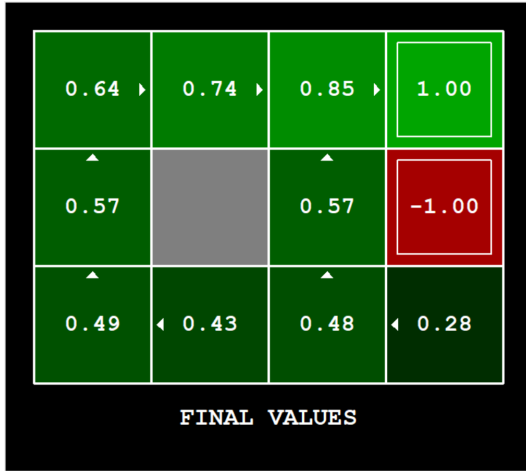The policy iteration is implemented in `policyIterationAgent.py`.

The results given by policy iteration under $\varepsilon = 0.01$ or $0.001$ are shown in Fig. 2.
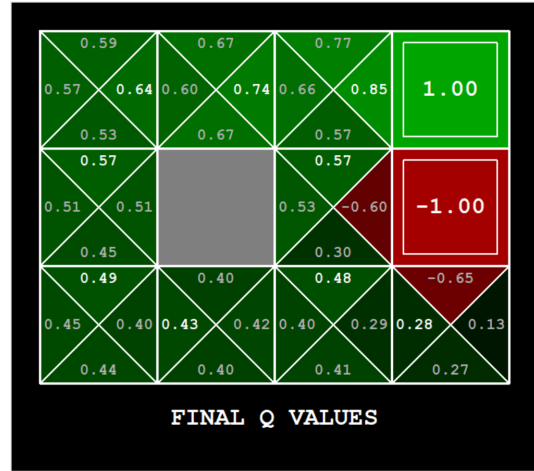
(a) Final Values when $\varepsilon = 0.01$
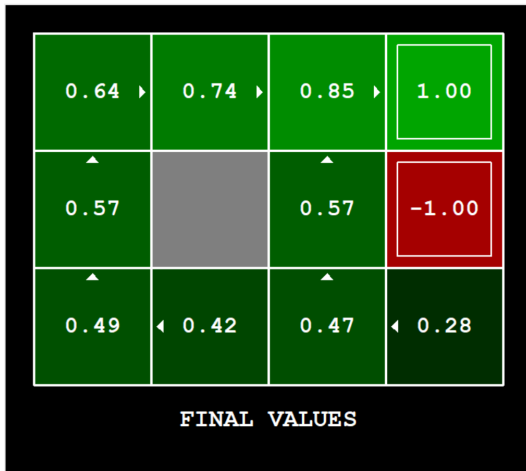


(b) Q-values when $\varepsilon = 0.01$
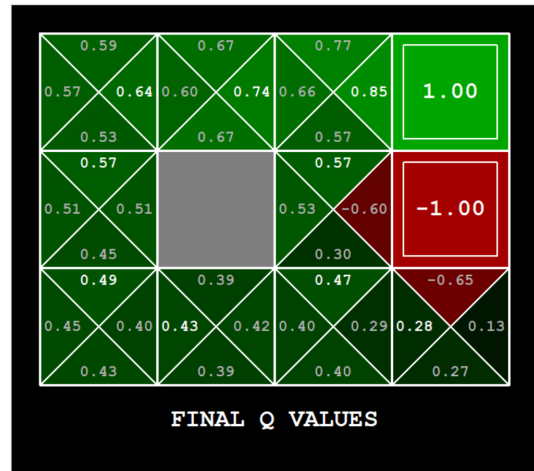


(c) Final Values when $\varepsilon = 0.001$



(d) Q-values when $\varepsilon = 0.01$

Figure 1: Value Iteration



(a) Final Values when $\varepsilon = 0.01$ or $0.001$



(b) Q-values when $\varepsilon = 0.01$ or $0.001$

Figure 2: Policy Iteration

## 2.3 Comparison of Speed of Convergence

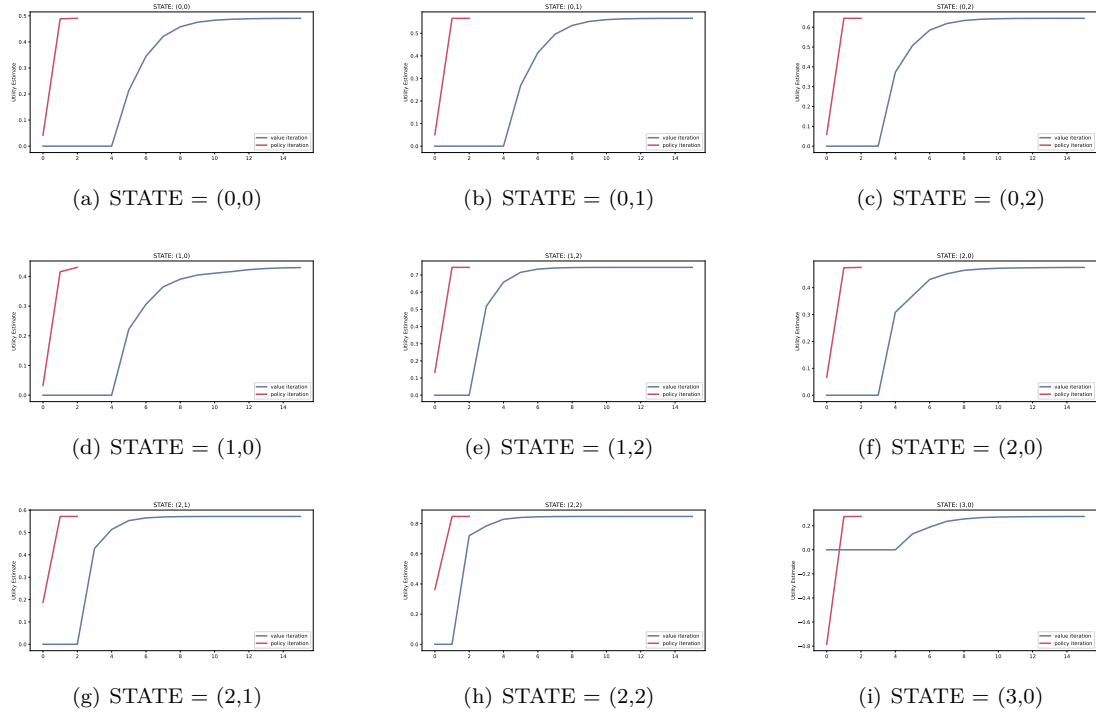The maximal change of value function during the iteration process is depicted in Fig. 3.

(a) STATE = (0,0)  (b) STATE = (0,1)  (c) STATE = (0,2)

(d) STATE = (1,0)  (e) STATE = (1,2)  (f) STATE = (2,0)

(g) STATE = (2,1)  (h) STATE = (2,2)  (i) STATE = (3,0)

Figure 3: Convergence of utilities of all states under value iteration and policy iteration

As shown above, in the aspect of utility, policy iteration **converges faster than** value iteration. Moreover, closer a state is to the terminal state, faster the value iteration of the utility of it converges.

\* Still, it seems a little unfair to compare in such ways. In the process of policy iteration, value estimate also iterates for several times to estimate the value function of the current policy. But it is also possible to estimate the value by solving a linear function, which won't lead to large time consumption.
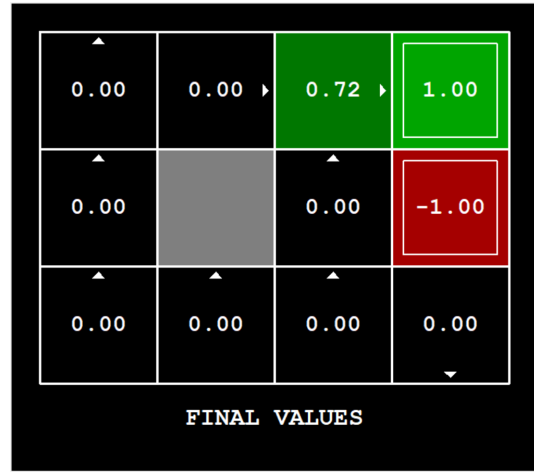
The policy actions of all states against the number of iterations are as follows. For policy iteration, policy actions are shown in Fig. 4. For value iteration, policy actions are shown in Fig. 5 and Fig. 6.
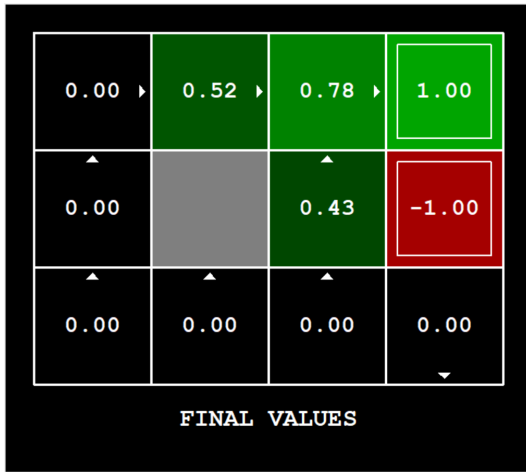
(a) iteration=0  (b) iteration=1  (c) iteration=2

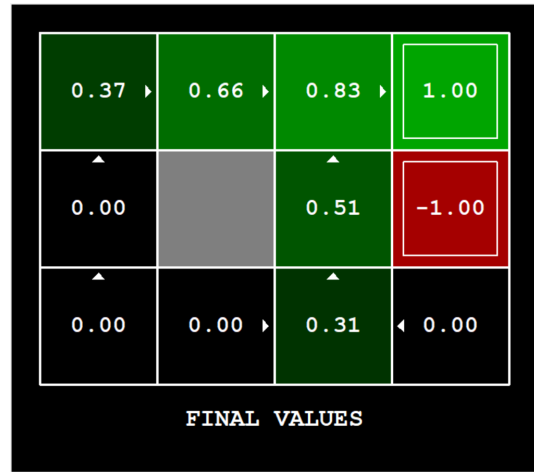Figure 4: Policy Actions given by Policy Iterations
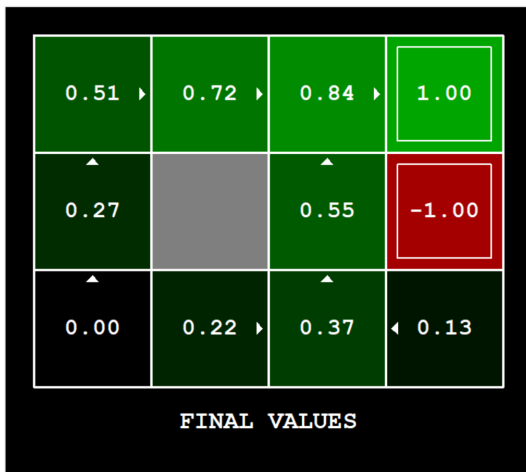
4

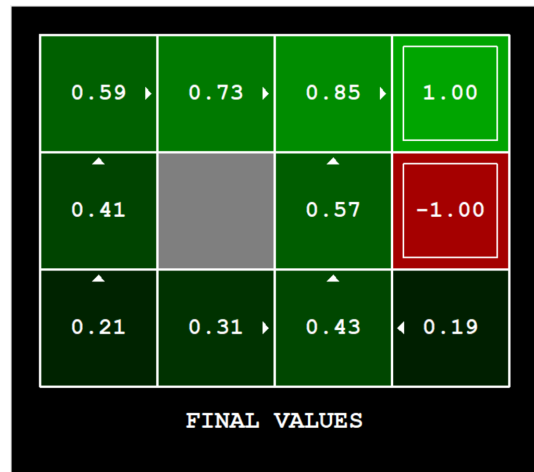(a) iteration = 0

(b) iteration = 1
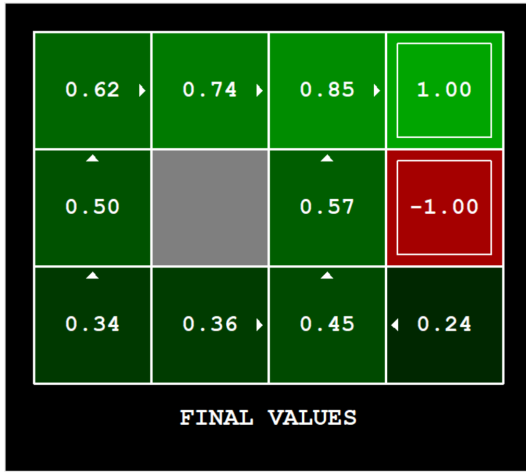
(c) iteration = 2

(d) iteration = 3
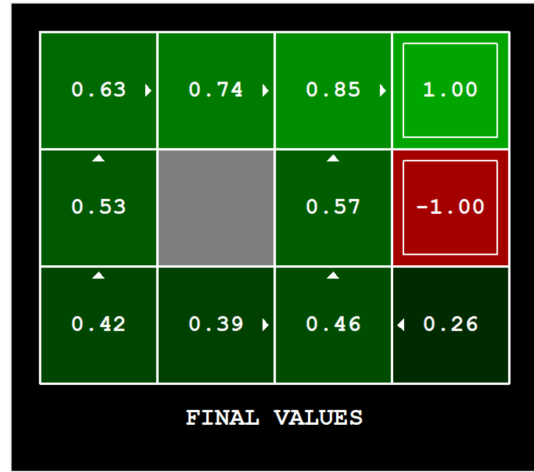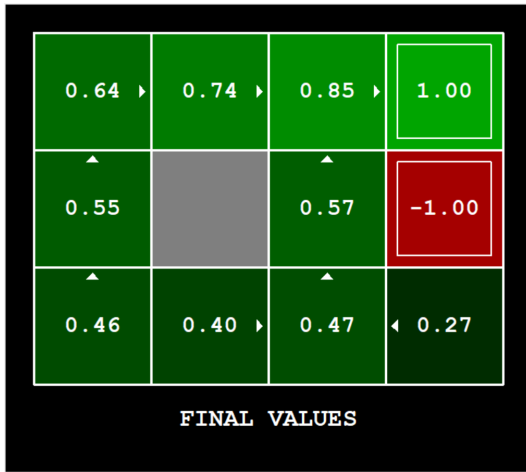
(e) iteration = 4

(f) iteration = 5

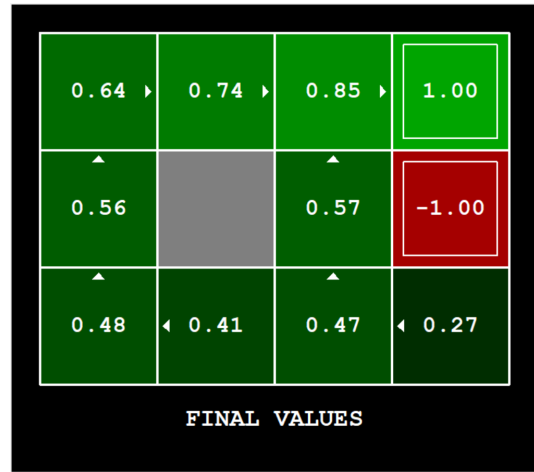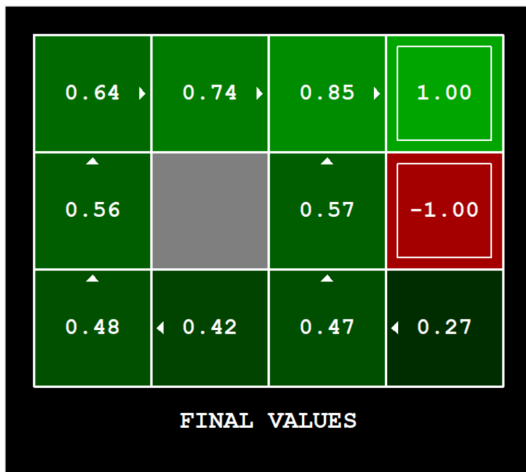Figure 5: Policy Actions given by Value Iterations (Part. I)
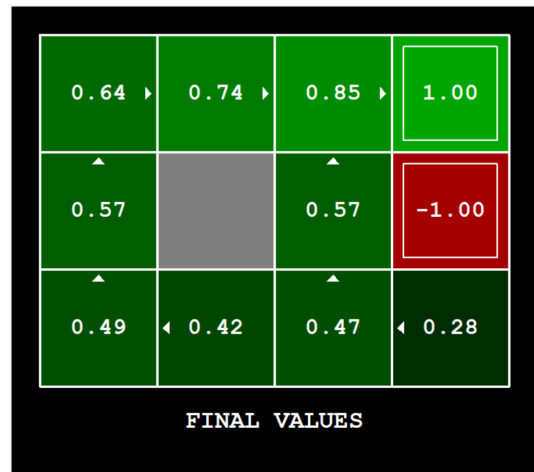
(a) iteration = 6

(b) iteration = 7

(c) iteration = 8

(d) iteration = 9

(e) iteration = 10

(f) iteration = 11

Figure 6: Policy Actions given by Value Iterations (Part. II)

(a) iteration = 12 & 13
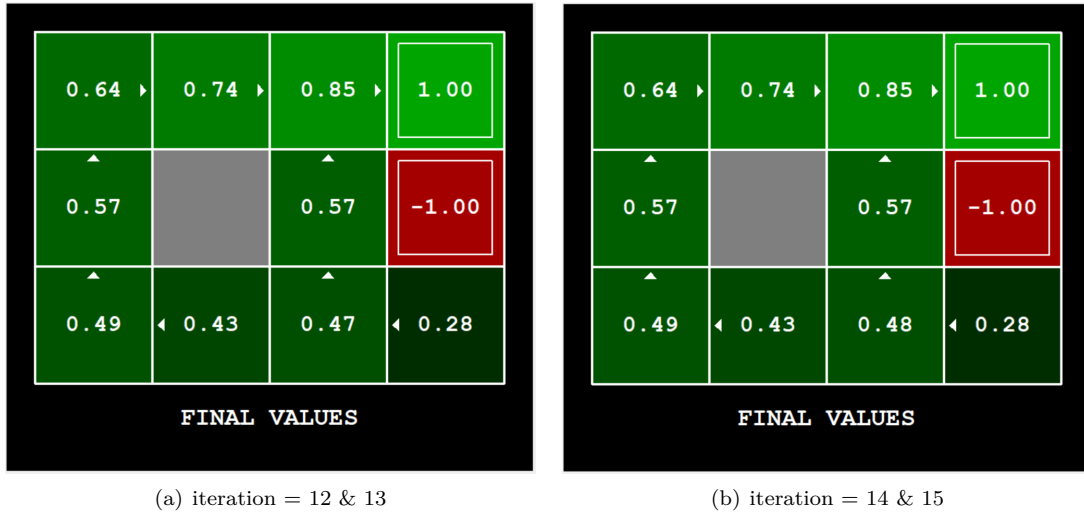
(b) iteration = 14 & 15

Figure 7: Policy Actions given by Value Iterations (Part. III)

The policy is stable and optimal after the 2nd policy iteration. Meanwhile, even the policy given by the 1st iteration is quite close to the final optimal policy. On the other hand, the policy converges to the optimal policy in the 9th value iteration. In fact, the 10th to 15th value iteration is trying to approximate the value function, though the policy is already stable and optimal.

Thus, value iteration **converges slower than** policy iteration.