



[Course](#) > [Week 6](#) > [Project...](#) > [p3\\_rl\\_q...](#)

## p3\_rl\_q1\_value\_iteration

### Question 1 (6 points): Value Iteration

Write a value iteration agent in `ValueIterationAgent`, which has been partially specified for you in `valueIterationAgents.py`. Your value iteration agent is an offline planner, not a reinforcement learning agent, and so the relevant training option is the number of iterations of value iteration it should run (option `-i`) in its initial planning phase.

`ValueIterationAgent` takes an MDP on construction and runs value iteration for the specified number of iterations before the constructor returns.

Value iteration computes  $k$ -step estimates of the optimal values,  $V_k$ . In addition to running value iteration, implement the following methods for `ValueIterationAgent` using  $V_k$ .

- `computeActionFromValues(state)` computes the best action according to the value function given by `self.values`.
- `computeQValueFromValues(state, action)` returns the Q-value of the (state, action) pair given by the value function given by `self.values`.

These quantities are all displayed in the GUI: values are numbers in squares, Q-values are numbers in square quarters, and policies are arrows out from each square.

*Important:* Use the "batch" version of value iteration where each vector  $V_k$  is computed from a fixed vector  $V_{k-1}$  (like in lecture), not the "online" version where one single weight vector is updated in place. This means that when a state's value is updated in iteration  $k$  based on the values of its successor states, the successor state values used in the value update computation should be those from iteration  $k-1$  (even if some of the successor states had already been updated in iteration  $k$ ). The difference is discussed in [Sutton & Barto](#) in the 6th paragraph of chapter 4.1.

*Note:* A policy synthesized from values of depth  $k$  (which reflect the next  $k$  rewards) will actually reflect the next  $k+1$  rewards (i.e. you return  $\pi_{k+1}$ ). Similarly, the  $Q$ -values will also reflect one more reward than the values (i.e. you return  $Q_{k+1}$ ).

You should return the synthesized policy  $\pi_{k+1}$ .

*Hint:* Use the `util.Counter` class in `util.py`, which is a dictionary with a default value of zero. Methods such as `totalCount` should simplify your code. However, be careful with `argMax`: the actual `argmax` you want may be a key not in the counter!

*Note:* Make sure to handle the case when a state has no available actions in an MDP (think about what this means for future rewards).

To test your implementation, run the autograder:

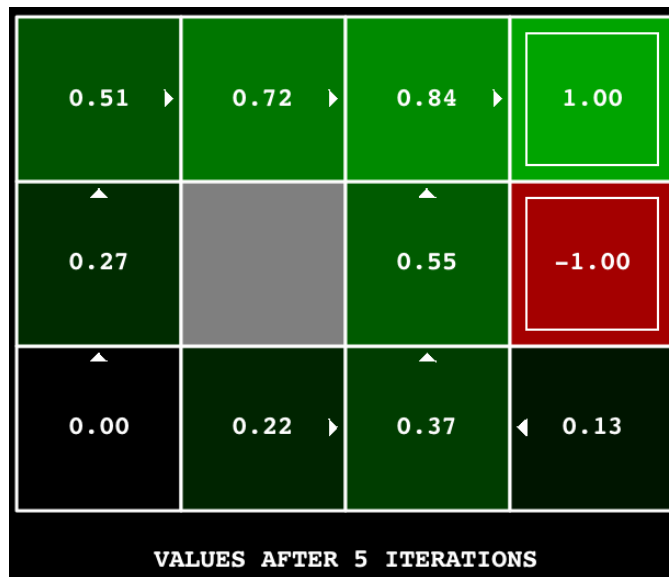
```
python autograder.py -q q1
```

The following command loads your `ValueIterationAgent`, which will compute a policy and execute it 10 times. Press a key to cycle through values,  $Q$ -values, and the simulation. You should find that the value of the start state ( $V(\text{start})$ ), which you can read off of the GUI) and the empirical resulting average reward (printed after the 10 rounds of execution finish) are quite close.

```
python gridworld.py -a value -i 100 -k 10
```

*Hint:* On the default `BookGrid`, running value iteration for 5 iterations should give you this output:

```
python gridworld.py -a value -i 5
```



*Grading:* Your value iteration agent will be graded on a new grid. We will check your values, Q-values, and policies after fixed numbers of iterations and at convergence (e.g. after 100 iterations).

© All Rights Reserved