

# Policy Iteration

Quiz, 7 questions

1  
point

1.

What are the two main steps in value-based approach to Reinforcement Learning?

- ☐ 2 - extract a reward function from the value function.
  - ☐ 1 - estimate a reward function.
  - ☐ 2 - extract a value function from the reward function.
  - ☒ 1 - build a value function.
  - ☐ 1 - build a policy function.
  - ☒ 2 - extract a policy function from the value function.
  - ☐ 2 - extract a value function from policy.
- 

1  
point

2.

What is true about policy improvement? Recall that,

total return = immediate reward + the discounted expected return from the next state under policy  $\pi$ .

- ☐ Making several policy improvements in a row may increase the performance of a new policy.
- ☐ An agent acts greedily with respect to the immediate reward only and ignores the remaining expected return under policy  $\pi$ .
- ☒ An agent acts greedily with respect to combination of immediate reward and the expected return under policy  $\pi$ .
- ☐ Relying on the estimates of expected return under policy  $\pi$  may lead to deterioration of an agent's performance in some states. This is so because the estimates will no longer valid as soon as policy is changed (improved).

---

# Policy Iteration

Quiz, 7 questions  
point

3.

How many different value functions can correspond to any particular policy function?

- ☐ Depends on number of actions.
- ☒ One
- ☐ Infinite
- ☐ Depends on number of states.

---

1  
point

4.

Why we don't need the precise solution of a system of Bellman equations?

- ☐ The solution of such system of equations is intractable on any modern supercomputer. Thus we have to approximate.
- ☐ The system of Bellman equations may have no solution at all. Thus we should be satisfied with an approximation.
- ☐ We want to sacrifice the global optimality for much faster convergence.
- ☒ After reaching some precision level further refinements of the solution will not change the result of subsequent policy improvement.

---

2  
points

5.

Generalised Policy Iteration (GPI)

- ☒ does not require to perform policy evaluation until its convergence
- ☐ requires to perform policy evaluation until convergence at every iteration.
- ☐ converges to local optimum.
- ☐ depends on initialization.

## Policy Iteration

Quiz, 7 questions

- ☒ does not require to improve policy in each and every state as long as policy in any state is improved once in a while
  - ☐ requires to improve policy in each and every state before subsequent policy evaluation.
  - ☒ does not depend on initialization.
- 

1  
point

6.

How can we recover the optimal policy solely from  $q^*$  function?

- ☐ With max operator.
  - ☐ Sample from a distribution that is proportional to  $q$ -values.
  - ☐ Find the action that is closest in  $q$ -value to average  $q$ -value over actions.
  - ☐ It is impossible without the knowledge of environment dynamics.
  - ☒ With argmax operator.
- 

1  
point

7.

What is the difference between Policy Iteration and Value Iteration?

- ☒ Policy Iteration updates value function until numerical convergence of all its state values before each policy improvement step.
  - ☐ Value Iteration updates value function until numerical convergence of all its state values before each policy improvement step.
  - ☐ Policy Iteration perform only one iteration of policy evaluation before policy improvement step.
  - ☒ Value Iteration perform only one iteration of policy evaluation before policy improvement step.
- 



I, **Jiadao Zhao**, understand that submitting work that isn't my own may result in permanent failure of this course or deactivation of my Coursera account.

Quiz, 7 questions

Learn more about Coursera's Honor Code

Submit Quiz

