

**Congratulations! You passed!**

TO PASS 80% or higher

Keep Learning

GRADE  
100%

## Value Functions and Bellman Equations

TOTAL POINTS 10

1. A policy is a function which maps \_\_\_ to \_\_\_.

1 / 1 point

- Actions to probability distributions over values.
- States to actions.
- States to probability distributions over actions.
- Actions to probabilities.
- States to values.

**Correct**

Correct!

2. The term "backup" most closely resembles the term \_\_\_ in meaning.

1 / 1 point

- Value
- Update
- Diagram

Correct

3. At least one deterministic optimal policy exists in every Markov decision process.

1 / 1 point

- False
- True

**Correct**

Correct! Let's say there is a policy  $\pi_1$  which does well in some states, while policy  $\pi_2$  does well in others. We could combine these policies into a third policy  $\pi_3$ , which always chooses actions according to whichever of policy  $\pi_1$  and  $\pi_2$  has the highest value in the current state.  $\pi_3$  will necessarily have a value greater than or equal to both  $\pi_1$  and  $\pi_2$  in every state! So we will never have a situation where doing well in one state requires sacrificing value in another. Because of this, there always exists some policy which is best in every state. This is of course only an informal argument, but there is in fact a rigorous proof showing that there must always exist at least one optimal deterministic policy.

4. The optimal state-value function:

1 / 1 point

- Is unique in every finite Markov decision process.
- Is not guaranteed to be unique, even in finite Markov decision processes.

**Correct**

Correct! The Bellman optimality equation is actually a system of equations, one for each state, so if there are N states, then there are N equations in N unknowns. If the dynamics of the environment are known, then in principle one can solve this system of equations for the optimal value function using any one of a variety of methods for solving systems of nonlinear equations. All optimal policies share the same optimal state-value function.

- Yes, adding a constant to all rewards changes the set of optimal policies.
- No, as long as the relative differences between rewards remain the same, the set of optimal policies is the same.

**Correct**

Correct! Adding a constant to the reward signal can make longer episodes more or less advantageous (depending on whether the constant is positive or negative).

6. Does adding a constant to all rewards change the set of optimal policies in continuing tasks?

1 / 1 point

- No, as long as the relative differences between rewards remain the same, the set of optimal policies is the same.
- Yes, adding a constant to all rewards changes the set of optimal policies.

**Correct**

Correct! Since the task is continuing, the agent will accumulate the same amount of extra reward independent of its behavior.

7. Select the equation that correctly relates  $v_*$  to  $q_*$ . Assume  $\pi$  is the uniform random policy.

1 / 1 point

- $v_*(s) = \sum_{a,r,s'} \pi(a|s)p(s', r|s, a)q_*(s')$
- $v_*(s) = \max_a q_*(s, a)$
- $v_*(s) = \sum_{a,r,s'} \pi(a|s)p(s', r|s, a)[r + q_*(s')]$
- $v_*(s) = \sum_{a,r,s'} \pi(a|s)p(s', r|s, a)[r + \gamma q_*(s')]$

**Correct**

Correct!

8. Select the equation that correctly relates  $q_*$  to  $v_*$  using four-argument function  $p$ .

1 / 1 point

- $q_*(s, a) = \sum_{s', r} p(s', r|s, a)[r + v_*(s')]$
- $q_*(s, a) = \sum_{s', r} p(s', r|s, a)\gamma[r + v_*(s')]$
- $q_*(s, a) = \sum_{s', r} p(s', r|s, a)[r + \gamma v_*(s')]$

**Correct**

Correct!

9. Write a policy  $\pi_*$  in terms of  $q_*$ .

1 / 1 point

- $\pi_*(s) = q_*(s, a)$
- $\pi_*(s) = \max_{a'} q_*(s, a')$
- $\pi_*(s) = 1 \text{ if } a = \arg\max_{a'} q_*(s, a'), \text{ else } 0$

**Correct**

Correct!

10. Give an equation for some  $\pi_*$  in terms of  $v_*$  and the four-argument  $p$ .

1 / 1 point

- $\pi_*(s) = 1 \text{ if } v_*(s) = \sum_{s', r} p(s', r|s, a)[r + \gamma v_*(s')], \text{ else } 0$
- $\pi_*(s) = \max_{a'} \sum_{s', r} p(s', r|s, a')[r + \gamma v_*(s')]$
- $\pi_*(s) = \sum_{s', r} p(s', r|s, a)[r + \gamma v_*(s')]$
- $\pi_*(s) = 1 \text{ if } v_*(s) = \max_{a'} \sum_{s', r} p(s', r|s, a')[r + \gamma v_*(s')], \text{ else } 0$