

1.

Question 1

The value of any state under an optimal policy is ____ the value of that state under a non-optimal policy. [Select all that apply]

Strictly greater than

Greater than or equal to

Correct! This follows from the policy improvement theorem.

Strictly less than

Less than or equal to

1 / 1 point

2.

Question 2

If a policy is greedy with respect to the value function for the

equiprobable random policy, then it is **guaranteed** to be an optimal policy.

True

False

Correct! Only policies greedy with respect to the optimal value function are guaranteed to be optimal.

1 / 1 point

3.

Question 3

Let

v_π

v

π

v , start subscript, π , end subscript

be the state-value function for the policy

π

π

π

. Let

π'

π

,

π , prime

be greedy with respect to

$\forall \pi$

ν

π

ν , start subscript, pi, end subscript

. Then

$$\nu\pi'\geq\nu\pi$$

ν

π

,

$$\geq\nu$$

π

$v_{\pi'}$ is greater than or equal to v_{π}

.

True

Correct! This is a consequence of the policy improvement theorem.

False

1 / 1 point

4.

Question 4

What is the relationship between value iteration and policy iteration? [Select all that apply]

Value iteration is a special case of policy iteration.

Value iteration and policy iteration are both special cases of

generalized policy iteration.

Correct!

Policy iteration is a special case of value iteration.

1 / 1 point

5.

Question 5

The word synchronous means "at the same time". The word asynchronous means "not at the same time". A dynamic programming algorithm is: [Select all that apply]

Asynchronous, if it does not update all states at each iteration.

Correct! Only algorithms that update every state exactly once at each iteration are synchronous.

Synchronous, if it systematically sweeps the entire state space at each iteration.

Correct! Only algorithms that update every state exactly once at each iteration are synchronous.

Asynchronous, if it updates some states more than others.

Correct! Only algorithms that update every state exactly once at each iteration are synchronous.

1 / 1 point

6.

Question 6

All Generalized Policy Iteration algorithms are synchronous.

True

False

Correct! A Generalized Policy Iteration algorithm can update states in a non-systematic fashion.

1 / 1 point

7.

Question 7

Which of the following is true?

Synchronous methods generally scale to large state spaces better than asynchronous methods.

Asynchronous methods generally scale to large state spaces better than synchronous methods.

Correct! Asynchronous methods can focus updates on more relevant states, and update less relevant states less often. If the state space is very large, asynchronous methods may still be able to achieve good performance whereas even just one synchronous sweep of the state space may be intractable.

1 / 1 point

8.

Question 8

Why are dynamic programming algorithms considered planning methods? [Select all that apply]

They use a model to improve the policy.

Correct! This is the definition of a planning method.

They compute optimal value functions.

They learn from trial and error interaction.

Incorrect. A planning method is a method that uses a model to improve the policy.

1 point

9.

Question 9

Consider the undiscounted, episodic MDP below. There are four actions possible in each state, $A = \{\text{up, down, right, left}\}$, which deterministically cause the corresponding state transitions, except that actions that would take the agent off the grid in fact leave the state unchanged. The right half of the figure shows the value of each state under the equiprobable random policy. If

π

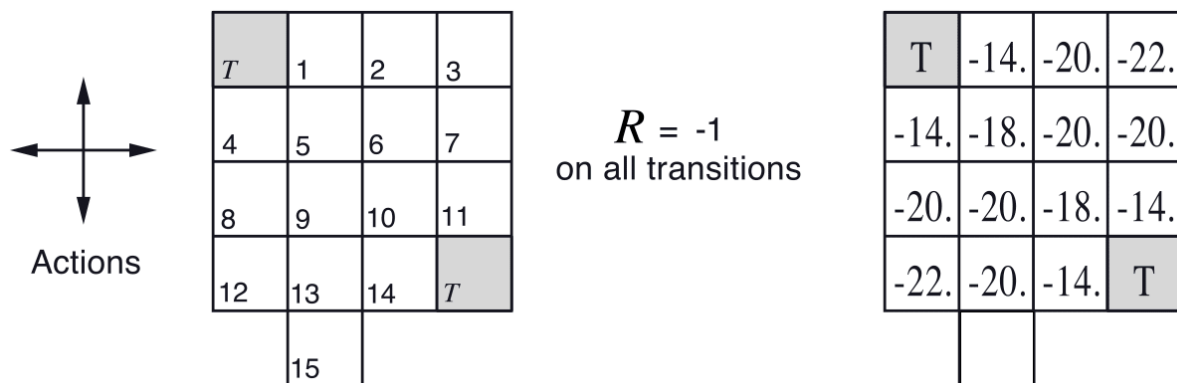
π

π

is the equiprobable random policy, what is

$q(7, \text{down})$

?



$q(7, \text{down}) = -14$

Incorrect. Moving down incurs a reward of -1 before reaching state 11, from which the expected future return is -14.

$q(7, \text{down}) = -20$

$$q(7, \text{down}) = -21$$

$$q(7, \text{down}) = -15$$

1 point

10.

Question 10

Consider the undiscounted, episodic MDP below. There are four actions possible in each state, $A = \{\text{up, down, right, left}\}$, which deterministically cause the corresponding state transitions, except that actions that would take the agent off the grid in fact leave the state unchanged. The right half of the figure shows the value of each state under the equiprobable random policy. If

π

π

π

is the equiprobable random policy, what is

$$v(15)$$

$$v(15)$$

v, left parenthesis, 15, right parenthesis

? Hint: Recall the Bellman equation

$$v(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v(s')]$$

$$v(s) = \sum$$

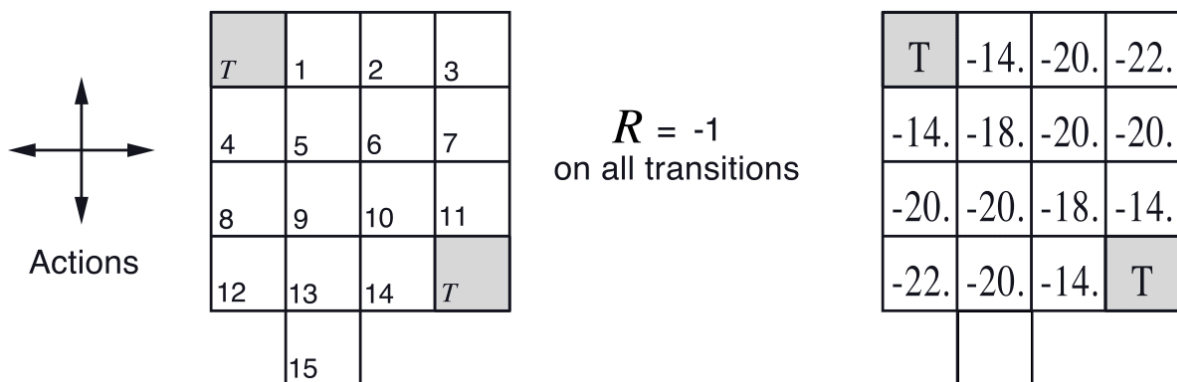
$$a$$

$$\pi(a|s) \sum$$

$$s',r$$

$$p(s', r | s, a) [r + \gamma v(s')]$$

v, left parenthesis, s, right parenthesis, equals, sum, start subscript, a, end subscript, pi, left parenthesis, a, vertical bar, s, right parenthesis, sum, start subscript, s, ', comma, r, end subscript, p, left parenthesis, s, ', comma, r, vertical bar, s, comma, a, right parenthesis, open bracket, r, plus, gamma, v, left parenthesis, s, ', right parenthesis, close bracket



$$v(15) = -25$$

$$v(15) = -25$$

v, left parenthesis, 15, right parenthesis, equals, minus, 25

$$v(15)=-21$$

$$v(15)=-21$$

v, left parenthesis, 15, right parenthesis, equals, minus, 21

$$v(15)=-23$$

$$v(15)=-23$$

v, left parenthesis, 15, right parenthesis, equals, minus, 23

$$v(15)=-24$$

$$v(15)=-24$$

v, left parenthesis, 15, right parenthesis, equals, minus, 24

Correct! We can get this by solving for the unknown variable

$v(15)$

$v(15)$

v, left parenthesis, 15, right parenthesis

. Let's call this unknown

x

x

x

. We solve for

x

x

x

in the equation

$$x = \frac{1}{4}(-21) + \frac{3}{4}(-1+x)$$

$$x = \frac{1}{4}(-21) + \frac{3}{4}(-1+x)$$

$x, equals, 1, slash, 4, left\ parenthesis, minus, 21, right\ parenthesis, plus, 3, slash,$
 $4, left\ parenthesis, minus, 1, plus, x, right\ parenthesis$

. The first term corresponds to transitioning to state

13

13

13

. The second term corresponds to taking one of the other three actions, incurring
a reward of

-1

-1

minus, 1

and staying in state

x

x

x

.

$$v(15)=-22$$

$$v(15)=-22$$

v, left parenthesis, 15, right parenthesis, equals, minus, 22

1 / 1 point

[Like](#)

[Dislike](#)

[Report an issue](#)