

Unit 5 Reinforcement Learning (2

Course > weeks)

> Homework 6 > 2. Q-Value Iteration

Audit Access Expires May 11, 2020

You lose all access to this course, including your progress, on May 11, 2020.

2. Q-Value Iteration

Consider an Markov Decision Process with 6 states $s\in\{0,1,2,3,4,5\}$ and 2 actions $a\in\{C,M\}$, defined by the following transition probability functions For states 1, 2, and 3:

$$T\left(s,M,s-1\right) =1$$

$$T\left(s,C,s+2\right)=0.7$$

$$T\left(s,C,s
ight)=0.3$$

For state 0:

$$T\left(s,M,s\right) =1$$

$$T(s, C, s) = 1$$

For states 4 and 5:

$$T\left(s,M,s-1\right)=1$$

$$T(s, C, s) = 1$$

Note that all transition probabilities not defined by the above are equal to 0. The rewards R are defined by:

$$R\left(s,a,s'
ight)=|s'-s|^{rac{1}{3}}\;orall s
eq s',$$
 and $R\left(s,a,s
ight)=(s+4)^{rac{-1}{2}}$, $orall s
eq 0.$ $R\left(0,M,0
ight)=R\left(0,C,0
ight)=0.$ Also, the discount factor $\gamma=0.6.$

We initialize $Q_{0}\left(s,a
ight) =0\ orall s\in \{0,1,2,3,4,5\}$ and $orall a\in \{C,M\}.$

1

1 point possible (graded)

We can conclude from this information that $\boldsymbol{0}$ is a terminal state.

True 🗸

False

Solution:

From the transition probabilities, we can see that no matter which action you take, once you are in state 0, you can never leave.

Submit

You have used 0 of 1 attempt

1 Answers are displayed within the problem

2

0.0/6.0 points (graded)

Input the Q-values $Q_{1}\left(s,a\right)$ correct to 3 decimal places after one Q-value iteration

 $Q_{1}\left(0,M
ight) =% {\displaystyle\int\limits_{0}^{\infty}} \left[\left(0,M
ight) -\left(\left(0,M
ight)
ight] \left(\left(0,M
ight)
ight) -\left(\left(\left(0,M
ight)
ight)
ight] \left(\left(\left(0,M
ight)
ight)
ight]$ Answer: 0

 $Q_{1}\left(0,C\right) =$ Answer: 0

 $Q_{1}\left(1,M
ight) = oxed{Q_{1}\left(1,C
ight) =}% oxed{Q_{1}\left(1,C
ight) =}% oxed{Q_{1}\left(1,C
ight) =}% oxed{Q_{2}\left(1,C
ight) =}% oxed{Q_{1}\left(1,C
ight) =}% oxed{Q_{2}\left(1,C
ight) =}% oxed{Q_$ Answer: 1

Answer: 1.016

$$Q_{1}\left(2,M\right) =$$

Answer: 1

$$Q_{1}\left(2,C\right) =%$$

Answer: 1.004

$$Q_{1}\left(3,M\right) =$$

Answer: 1

$$Q_{1}\left(3,C\right) =%$$

Answer: 0.995

$$Q_{1}\left(4,M\right) =$$

Answer: 1

$$Q_{1}\left(4,C\right) =%$$

Answer: 0.354

$$Q_{1}\left(5,M\right) =%$$

Answer: 1

$$Q_{1}\left(5,C\right) =$$

Answer: 0.333

Solution:

1.
$$Q_{1}\left(0,M
ight)$$
: $Q_{1}\left(0,M
ight)=0$ because $R\left(0,M,0
ight)=0$ and $T\left(0,M,s'
ight)=0\ orall s'
eq0$

2.
$$Q_{1}\left(0,C
ight)$$
: $Q_{1}\left(0,C
ight)=0$ because $R\left(0,C,0
ight)=0$ and $T\left(0,C,s'
ight)=0$ $orall s'
eq 0$

3.
$$Q_1\left(1,M
ight):\left|\left(0-1
ight)^{rac{1}{3}}
ight|=1$$

4.
$$Q_1\left(1,C\right)$$
: $0.7*\left|\left(3-1\right)^{rac{1}{3}}\right|+0.3*5^{rac{-1}{2}}=0.882+0.134=1.016$

5.
$$Q_{1}\left(2,M\right)$$
: Just as in $Q_{1}\left(1,M\right)$

Answer: 0

6.
$$Q_1\left(2,C
ight)$$
: $0.7*\left|\left(3-1
ight)^{rac{1}{3}}
ight|+0.3*5^{rac{-1}{2}}=0.882+0.122=1.004$

- 7. $Q_{1}\left(3,M\right)$: Just as in $Q_{1}\left(1,M\right)$
- 8. $Q_1\left(3,C
 ight)$: $0.7*\left|\left(3-1
 ight)^{rac{1}{3}}
 ight|+0.3*5^{rac{-1}{2}}=0.882+0.113=0.995$
- 9. $Q_{1}\left(4,M
 ight)$: Just as in $Q_{1}\left(1,M
 ight)$
- 10. $Q_1\left(4,C
 ight)$: $8^{rac{-1}{2}}=0.354$
- 11. $Q_{1}\left(5,M\right) :$ Just as in $Q_{1}\left(1,M\right)$
- 12. $Q_1\left(5,C\right)$: $9^{rac{-1}{2}}=0.333$

Submit

You have used 0 of 4 attempts

1 Answers are displayed within the problem

3

0.0/3.0 points (graded)

What are the values $V_{1}\left(s\right)$ corresponding to $Q_{1}\left(s,a\right)$?

$$V_{1}\left(0\right) =$$

$$V_{1}\left(1
ight)=% \left(1-\frac{1}{2}\right)^{2}$$
 Answer: 1.016

$$V_{1}\left(2
ight) =% \left\{ V_{1}\left(2
ight) =\left[1.004
ight] \right] \left\{ 1.004
ight\} \left[1.004
ight] \left[1.00$$

$$V_{1}\left(3
ight) =% \left\{ \left\{ V_{1}\left(3
ight) \right\} \right\} \left\{ \left\{ V_{1}\left(3
ight) \right\} \left\{ V_{1}\left(3
ight) \right\} \left\{ V_{1}\left(3
ight) \right\} \left\{ V_{1}\left(3
ight) \right\} \left\{ V_{1}\left(3
ight) \right\} \left\{ V_{1}\left(3
ight) \right\} \left\{ V_{1}\left(3
ight) \right\} \left\{ \left\{ V_{1}\left(3
ight) \left\{ V_{1}\left(3
ight) \right\} \left\{ V_{1}\left(3
ight)$$

 $V_{1}\left(4
ight) =% {\displaystyle\int\limits_{0}^{\infty }} \left[{\left| {{
m V}_{1}}
ight|}
ight] \left[{\left| {{
m V}_$

Answer: 1

$$V_1\left(5
ight) =$$

Answer: 1

Solution:

Because: $V_{1}\left(s
ight)=\max_{a}Q_{1}\left(s,a
ight)$

Submit

You have used 0 of 2 attempts

1 Answers are displayed within the problem

4

5 points possible (graded)

What are the optimal policies we get from $Q_1\left(s,a\right)$?

 $\pi^*(1) =$



 \bigcirc M

 $\pi^*(2) =$



 \bigcirc M

$$\pi^*(3) =$$

1	-	
		1



$$\pi^*(4) =$$





$$\pi^*(5) =$$





Solution:

We pick the policy corresponding to the $V_{1}\left(s
ight)$ i.e. $\pi^{*}(s)$ = $\mathop{argmax}\limits_{a}Q_{1}\left(s,a
ight)$

Submit

You have used 0 of 2 attempts

1 Answers are displayed within the problem

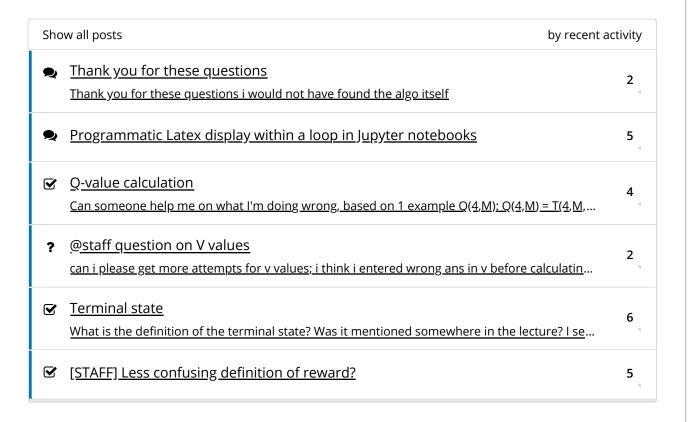
Discussion

Hide Discussion

Topic: Unit 5 Reinforcement Learning (2 weeks) :Homework 6 /

2. Q-Value Iteration

Add a Post



© All Rights Reserved

8 of 8 2020-05-09, 9:53 a.m.