

Report of Project3-Blackjack

Ziqin Luo 18307130198

November 26, 2020

1. Problem 1 Value Iteration

- a) Under the settings of problem 1, the formula of Value Iteration should be in the form of the format below:

$$V_{opt}(s) = \max_a \{R(s, a) + \gamma Q(s, a)\} = \max_a \{R(s, a) + \gamma \sum_a P(s'|s, a) V_{opt}(s')\}$$

To calculate the whole iterations, first we need to derive the **Reward Model** and the **Transition Model**. According to the conditions provided in the description, we have:

Reward Model:

R(s, a)	s = -1	s = 0	s = 1
a = -1	15	-5	16
a = +1	12.5	-5	26.5

Transition Model:

a = -1	s' = -2	s' = -1	s' = 0	s' = 1	s' = 2
s = -2					
s = -1	0.8		0.2		
s = 0		0.8		0.2	
s = 1			0.8		0.2
s = 2					

a = +1	s' = -2	s' = -1	s' = 0	s' = 1	s' = 2
s = -2					
s = -1	0.7		0.3		
s = 0		0.7		0.3	
s = 1			0.7		0.3
s = 2					

Besides, we know that $\gamma = 1$. Thus, we can calculate the Value Iteration in the form of a table:

Iteration	s = -2	s = -1	s = 0	s = +1	s = +2
-----------	--------	--------	-------	--------	--------

0	0	0	0	0	0
a = -1		15	-5	16	
a = +1		12.5	-5	26.5	
opt_action		a = -1	a = -1	a = +1	
1	0	15	-5	26.5	0
a = -1		14	12.3	12	
a = +1		11	13.45	23	
opt_action		a = -1	a = +1	a = +1	
2	0	14	13.45	23	0

Note: we always take the Values of terminal states as 0

From the table above, we know that after 0 iteration, which is equivalent to initialization, the values of each state:

$$V_{opt}(-2) = V_{opt}(2) = V_{opt}(-1) = V_{opt}(0) = V_{opt}(1) = 0$$

After the 1st iteration:

$$V_{opt}(-2) = V_{opt}(2) = 0, V_{opt}(-1) = 15, V_{opt}(0) = -5, V_{opt}(1) = 26.5$$

After the 2nd iteration:

$$V_{opt}(-2) = V_{opt}(2) = 0, V_{opt}(-1) = 14, V_{opt}(0) = 13.45, V_{opt}(1) = 23$$

b) We can derive the answers directly from question (a) :

$$\pi_{opt}(-1) = -1, \quad \pi_{opt}(0) = +1, \quad \pi_{opt}(1) = +1$$

2. Problem 2 Transforming MDPs

- a) Answer: The inequality $V_1(s_{start}) \geq V_2(s_{start})$ can't always holds. The counter example I constructed in [submission.py](#) is just the same as the simple MDP in Problem 1. Please refer to the source code for more details.
- b) We know that the MDP we have is acyclic. It implies that for any state s , it will only appear in the Search Trees of our MDP problem at some specific layer for once. That's to say we can store all the (s, a, s') triples in a list container and sort them in a topological order so that the (s_1, a, s'_1) triple with a deeper state s_1 (away from the root) comes before the (s_2, a, s'_2) triple with a shallower state s_2 (close to the root). Then, we can visit these triples successively and calculate their V_{opt} . This is equivalent to the operations that we start calculating the V_{opt} of states at the bottom of the MDP Search Tree and then calculate the V_{opt} of their father nodes, until we calculate the V_{opt} of root node in the Search Tree. It is obvious that to update the V_{opt} of all the states, we just need only a single pass over all the (s, a, s') triples.

c) We know that:

$$\mathbf{States}' = \mathbf{States} \cup \{\mathbf{o}\}, \mathbf{Actions}'(s) = \mathbf{Actions}(s), \gamma' = 1, \gamma < 1$$

For $\forall s \in \mathbf{States}$, we define set

$$S'_s = \{s' : T(s, a, s') > 0, s \text{ is NOT a terminal state}, a \in \mathbf{Actions}(s)\}.$$

$$S'_{s_{new}} = S'_s \cup \{\mathbf{o}\}$$

That is to say we see state \mathbf{o} as a possible successor state of all states in and we treat it as a terminal state, which implies that $V_{opt}(\mathbf{o}) = 0$.

Then, we redefine the new Reward Model $\mathbf{R}'(s, a, s')$ and Transition Model $\mathbf{T}'(s, a, s')$ as below:

- For primal model $\mathbf{T}(s, a, s')$ and new model $\mathbf{T}'(s, a, s')$:
 - Let $\mathbf{T}'(s, a, s') = \gamma \mathbf{T}(s, a, s')$.
 - $\forall s \in \mathbf{States}$ (s not terminal), let $\mathbf{T}'(s, a, \mathbf{o}) = 1 - \sum_{s' \neq \mathbf{o}} \mathbf{T}'(s, a, s') = 1 - \gamma$
- For primal model $\mathbf{R}(s, a, s')$ and new model $\mathbf{R}'(s, a, s')$:
 - Let $\mathbf{R}'(s, a, s') = \frac{1}{\gamma} \mathbf{R}(s, a, s')$
 - $\forall s \in \mathbf{States}$ (s not terminal), let $\mathbf{R}'(s, a, \mathbf{o}) = 0$

Then we have for $\forall s \in \mathbf{State}$ and s is not a terminal state:

$$\begin{aligned}
 V_{opt}(s) &= \max_{a \in \mathbf{Actions}(s)} \{R(s, a) + \gamma Q(s, a)\} \\
 &= \max_{a \in \mathbf{Actions}(s)} \left\{ R(s, a) + \gamma \sum_{s'} T(s, a, s') V_{opt}(s') \right\} \\
 &= \max_{a \in \mathbf{Actions}(s)} \left\{ \sum_{s'} \gamma T(s, a, s') \frac{1}{\gamma} R(s, a, s') + \sum_{s'} \gamma T(s, a, s') V_{opt}(s') \right\} \\
 &= \max_{a \in \mathbf{Actions}(s)} \left\{ \sum_{s' \in S'_s} T'(s, a, s') R'(s, a, s') + T'(s, a, \mathbf{o}) R'(s, a, \mathbf{o}) \right. \\
 &\quad \left. + \sum_{s' \in S'_s} T'(s, a, s') V_{opt}(s') + T'(s, a, \mathbf{o}) V_{opt}(\mathbf{o}) \right\} \\
 &= \max_{a \in \mathbf{Actions}'(s)} \left\{ \sum_{s' \in S'_{s_{new}}} T'(s, a, s') R'(s, a, s') + \sum_{s' \in S'_{s_{new}}} T(s, a, s') V_{opt}(s') \right\}
 \end{aligned}$$

$$= \max_{a \in \text{Actions}'(s)} \{R'(s, a) + \mathbf{1} * Q'(s, a)\} \quad (*)$$

From (*), we know the new MDP we build is equivalent to the original MDP.

3. Problem 4: Learning to Play Blackjack

Question (b) :

To illustrate the comparison result between Q-learning policy and the policy learned by Value Iteration, we need to finish the implementation of `simulate_QL_over_MDP()` in `submission.py`. Please refer to the source code for more details. There are detailed comments beside the codes.

After running the function `run4bHelper()` in `grader.py`, we get the match rate between between Q-learning policy and the policy learned by Value Iteration which is showed below:

```
----- START PART 4b-helper: Helper function to run Q-learning simulations for question 4b.
ValueIteration: 5 iterations
The match rate of using identityFeatureExtractor for smallMDP: 0.9737
Number of different actions: 1 Number of total actions: 38
ValueIteration: 15 iterations
The match rate of using identityFeatureExtractor for largeMDP: 0.7010
Number of different actions: 882 Number of total actions: 2950
----- END PART 4b-helper [took 0:00:03.405997 (max allowed 1000 seconds), 0/0 points]
```

We can see that when we run `simulate()` on `smallMDP`, the match rate is around 97% and only 1 state they produce a different action. We run `simulate()` on `largeMDP`, however, the performance gets worse and match rate is only around 70% and 882 states they produce a different action.

The reasons why the performance degrade, can be explained by the feature extractor we use. Refer to the implementation of `identityFeatureExtractor()`, we find that this feature extractor is just an indicator for the states in State Space of MDP. This implies that only when we encounter a transition process that contains the (s, a) state-action pair in the simulation, can we update the parameter $w_{s,a}$ which are useful to get $Q(s, a)$. When the complexity of MDP grows, the size of State Space grows and this leads to more new state-action pairs appear and to evaluate their Q values we need to simulate more transition processes that contain these new state-action pairs to update the parameters $w_{s,a}$. Therefore, when we run `simulate()` on `largeMDP` but still with only 30000 trials, there will be a lot more (s, a) pairs that may not be encountered in the trials or they appeared too few. Both of the two situations lead to the insufficient updates for the parameters. This accounts for the reasons why we can't have relatively accurate evaluations for their Q values

to get the optimal policy.

Question (d) :

To better answer question (d), we need to finish the implementation of `compare_changed_MDP()` in [submission.py](#). Please refer to the source code for more details. (For the reproductivity of this trial, I have set the random seed as 123.)

The screenshot below shows the results of this comparison trial:

```
----- START PART 4d-helper: Helper function to compare rewards when simulating RL over  
ValueIteration: 5 iterations  
*** Expected return for FixedRLAlgorithm (numTrials=30): 6.8333 > ***  
*** Expected return for FixedRLAlgorithm (numTrials=30000): 6.8309 >***  
*** Expected return for QLearningRLAlgorithm (numTrials=30): 0.7000 > ***  
*** Expected return for QLearningRLAlgorithm (numTrials=30000): 9.2527 > ***  
----- END PART 4d-helper [took 0:00:03.698559 (max allowed 60 seconds), 0/0 points]
```

As we can see, the expected reward with the policy provided by original MDP is **6.8309**(30 trials) and **6.8309**(30000 trials) while it is **0.7000**(30 trials) and **9.2527**(30000trials) when using Q-Learning.

From the results, we know two things:

- From [submission.py](#), we know that the threshold of `newThresholdMDP` is 15 but by using optimal policy of original MDP, the expected reward is far less than threshold. This implies that the optimal policy of original MDP is too **conservative** to play blackjack under the new settings.

This tells us that even though two MDPs maybe similar, the optimal policy of one MDP may behave badly in another one.

- Q-learning is run directly on the `newThresholdMDP`. When just do a few trials, Q-learning has too few samples to learn and its performance is very bad. However, if we conduct more trials and we find that Q-learning has improved dramatically and has a relatively better performance than the `fixedPolicy` does.

This implies that Q-learning can improve itself step by step using the samples, which is equivalent to find a better policy that is more suitable for `newThresholdMDP` step by step, and finally get a higher expected reward than the `fixedPolicy`.

Therefore, we know that Q-learning can become better and better through trials and perform better and better while the `fixedPolicy` always stays the same.