

EX 3

1) a) for stochastic case: $V_{\pi^*}^*(s) = \sum_a \pi(a|s) q_{\pi^*}^*(a, s)$

for deterministic case: $V_{\pi^*}^*(s) = \max_a q_{\pi^*}^*(a, s)$

b) $q^*(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma V^*(s')]$

c) $\pi^*(a|s) : \arg \max_a q_{\pi^*}^*(s, a)$

or $a^* = \arg \pi^*(a^*|s) = \arg \max_a (s, a)$

d) $\pi^*(a|s) : \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V^*(s')]$

or $a^* = \arg \pi^*(a^*|s) = \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V^*(s')]$

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

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

$q_{\pi}(s, a) = r(s, a) + \gamma \sum_{s'} p(s' | s, a) \sum_{a'} q_{\pi}(a', s') \pi(a' | s')$

$q_{\pi^*}^*(s, a) = r(s, a) + \gamma \sum_{s'} p(s' | s, a) \sum_{a'} \pi^*(a' | s') q_{\pi^*}^*(s', a')$

$p(s' | s, a) = \sum_r p(s', r | s, a) \quad r(s, a) = \sum_{s'} r p(s', r | s, a)$

Q2)

2) a) In policy improvement part, instead of comparing new policy and old policy for each state, we should compare value function for both new and old policies like below

3. Policy Improvement

Policy-stable \leftarrow true

For each $s \in \mathcal{S}$

$$V_{old}(s) \leftarrow V_{new}(s)$$

$$V_{new}(s) \leftarrow \arg \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma V_{old}(s')]$$

$$\text{If } \sum P(s', r | s, \pi_{new}(s)) [r + \gamma V_{new}(s')] >$$

$$\sum P(s', r | s, \pi_{old}(s)) [r + \gamma V_{old}(s')],$$

then Policy-stable \leftarrow false

If Policy-stable, then stop and return V^* and π^* ;

else go to policy evaluation part

In this way we never get stuck in an infinite loop

2) b) Yes It is possible to occur also for Value Iteration.
When in episodic tasks the rewards all are positive so $V(s)$ will always increase and agent never goes to terminal state so we need to consider negative reward in this case or use discount factor.

It is also more likely in non episodic task when there is no discount factor (exactly like Q4 - part C) so we need discount factor or generally we could use limit number of iteration.

Q3)

3.3) a)

1. Initialization

$Q(s,a) \leftarrow R \times R, \forall s \in S, \forall a \in A(s)$ arbitrarily
except $Q(s_{\text{terminal}}, a) = 0$

2. Policy Evaluation

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in S$

Loop for each $a \in A(s)$

$$q \leftarrow Q(s,a)$$

$$Q(s,a) \leftarrow \sum_{s',r} p(s',r|s,a) [r + \gamma \max_{a'} Q(s',a')]$$

$$\Delta \leftarrow \max(\Delta, |q - Q(s,a)|)$$

until $\Delta < \epsilon$

3. Policy Improvement

Policy-stable \leftarrow True

For each $s \in S$

old-action $\leftarrow \pi(s)$

$\pi(s) \leftarrow \arg \max_a Q(s,a)$

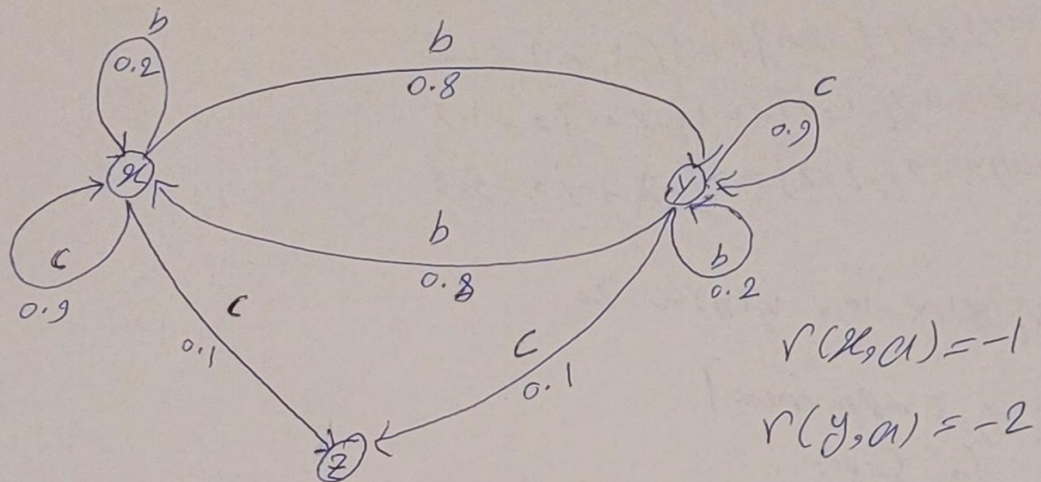
If old-action $\neq \pi(s)$, then Policy-stable \leftarrow False

If Policy-stable, then stop and return V^*, π^* ; else go to 2

Q 3-b and 4)

$$3) b) Q_{k+1}(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} Q_k(s', a')]$$

4) a)



In $S=X$, we never want to go to $S=Y$ in which for each action we will have -1 reward. In addition, we want to reach $S=Z$ which is terminal state as soon as possible. Therefore, in $S=X$ the best action is 'c'!

In $S=Y$, although the probability of reaching $S=Z$ is ~~the~~ same as $S=X$, the reward for each action is -2 , double comparing with reward in $S=X$. So, we prefer to go to $S=X$ and try to go to $S=Z$ in $S=X$ instead of $S=Y$ so the best action in $S=Y$ is 'b'!

4) b) Initialization:

$$\hat{V}(X)=0, \hat{V}(Y)=0, \hat{V}(Z)=0 \text{ (always)}, \gamma=1$$

Value evaluation:

$$\hat{V}_1(X) = 0.9[-1+0] + 0.1[-1+0] = -1$$

$$\hat{V}_1(Y) = 0.9[-2+0] + 0.1[-2+0] = -2$$

$$\hat{V}_2(X) = 0.9[-1-1] + 0.1[-2+0] = -1.2$$

$$\hat{V}_2(Y) = 0.9[-2-2] + 0.1[-2+0] = -3.8$$

⋮

$$\hat{V}_i(X) \sim -10, \hat{V}_i(Y) \sim -20$$

Policy Improvement:

In $S=X$:

$$\pi(X): c \Rightarrow \hat{V}(X) = 0.9[-1-10] + 0.1[-1+0] = -10$$

$$\pi(X): b \Rightarrow \hat{V}(X) = 0.8[-1-20] + 0.2[-1+10] = -19$$

In $S=Y$:

$$\pi(Y): c \Rightarrow \hat{V}(Y) = 0.9[-2-20] + 0.1[-2+0] = -20$$

$$\pi(Y): b \Rightarrow \hat{V}(Y) = 0.8[-2-10] + 0.2[-2+20] = -14$$

$$\pi(X) \leftarrow \underset{a}{\operatorname{argmax}} \hat{V}(X) = c$$

$$\pi(Y) \leftarrow \underset{a}{\operatorname{argmax}} \hat{V}(Y) = b$$

do again value evaluation till

$$\hat{V}(X) \sim -10, \hat{V}(Y) \sim -12.95$$

Policy Improvement

$S=X$

$$\pi(X): c \Rightarrow \hat{V}(X) = -10$$

$$\pi(X): b \Rightarrow \hat{V}(X) = -13.2$$

$S=Y$:

$$\pi(y): c \Rightarrow V(y) = -13.475$$

$$\pi(y): b \Rightarrow V(y) = -12.75$$

$$\pi(x) \leftarrow \arg \max_a V(x) = c$$

$$\pi(y) \leftarrow \arg \max_a V(y) = b$$

$$\pi_{old}(x) = \pi_{new}(x), \quad \pi_{old}(y) = \pi_{new}(y)$$

Stop and return $\pi^*(x): c, \pi^*(y): b$
 $V^*(x) = -10, V^*(y) = -12.75$

4) c) If the initial policy for both x, y states is b so we will have

$$V(x) = -1 + 0.2V(x) + 0.8V(y)$$

$$V(y) = -2 + 0.8V(x) + 0.2V(y)$$

which are inconsistent so there will be no solution for these equations.

with heavy discount these two equations will not more be inconsistent so we will find a unique answer. However, very small γ cause we have not optimal policy. For example in this question small γ will merge to 0 in 2 so cause we select action ' c ' even in $S=Y$, which is not optimal policy

For Questions 5 and 6:

- Both the gridworld and car rental environments are defined in `env.py`.
- The algorithms related to Q5 and Q6 are saved as `Q5_Gridworld.py` and `Q6_JacksCarRental.py`, you just need run these two py files

Q5)

A) Iterative policy evaluation

The value function associated with the random policy. As you can see it matches figure 3.2

```
random policy value function|
[[ 3.3  8.8  4.4  5.3  1.5]
 [ 1.5  3.   2.2  1.9  0.5]
 [ 0.   0.7  0.7  0.4 -0.4]
 [-1.  -0.4 -0.4 -0.6 -1.2]
 [-1.9 -1.3 -1.2 -1.4 -2. ]]
```

B) Value iteration

```
value iteration
0 = left, 1 = down, 2 = right, 3 = up
[[22.  24.4 22.  19.4 17.5]
 [19.8 22.  19.8 17.8 16. ]
 [17.8 19.8 17.8 16.  14.4]
 [16.  17.8 16.  14.4 13. ]
 [14.4 16.  14.4 13.  11.7]]
[[2. 0. 0. 0. 0.]
 [2. 3. 0. 0. 0.]
 [2. 3. 0. 0. 0.]
 [2. 3. 0. 0. 0.]
 [2. 3. 3. 3. 3.]]
```

C) Policy iteration

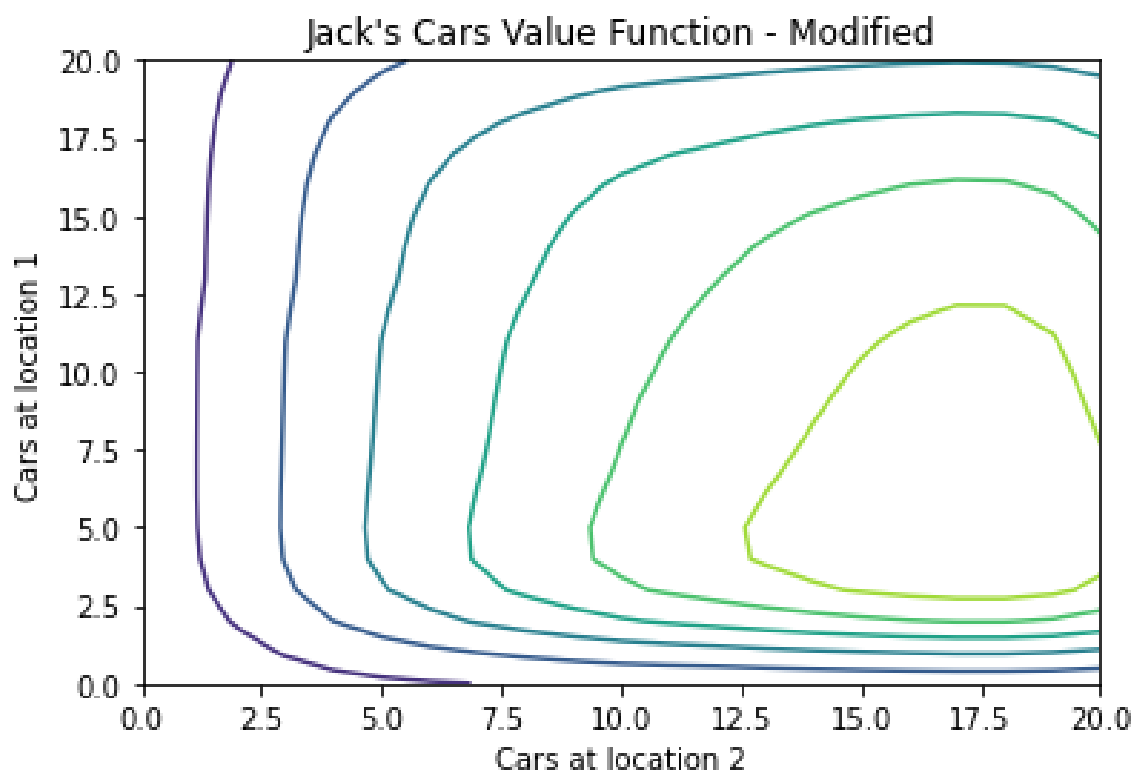
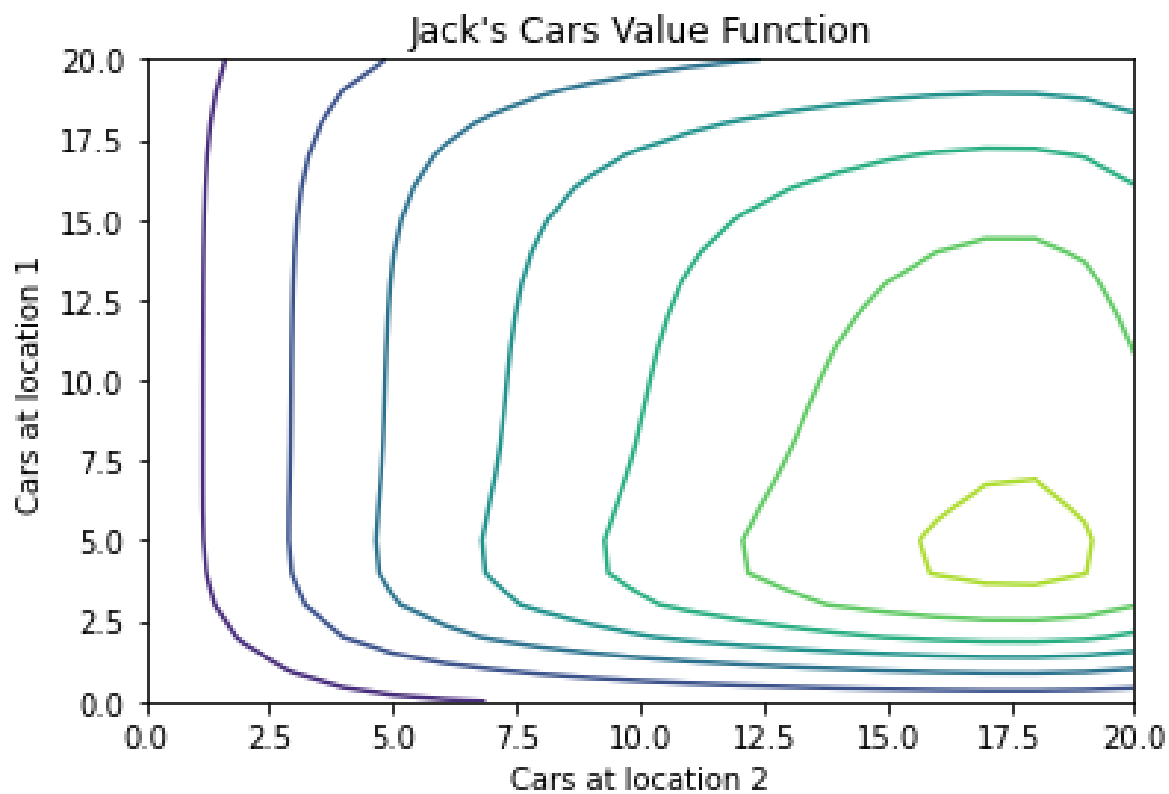
```
policy iteration, should look the same as value iteration
0 = left, 1 = down, 2 = right, 3 = up
[[22.  24.4 22.  19.4 17.5]
 [19.8 22.  19.8 17.8 16. ]
 [17.8 19.8 17.8 16.  14.4]
 [16.  17.8 16.  14.4 13. ]
 [14.4 16.  14.4 13.  11.7]]
[[2.  0.  0.  0.  0.]
 [2.  3.  0.  0.  0.]
 [2.  3.  0.  0.  0.]
 [2.  3.  0.  0.  0.]
 [2.  3.  3.  3.  3.]]
```

The optimal policy and the associated value function, which match figure 3.5

Q 6)

6. A. Unfortunately, because of an issue in my code, my policy converges to “0” for every state (no movement). So, there is no point at seeing a policy contour map. I provide the final value maps instead, which must also be incorrect.

B. Although I could not produce the correct policy, I could guess that because one piece of movement from lot 1 to 2 is free, the policy will be more willing to move when lot 2 is low. Additionally, we will see more movement in the 6-14 car midrange as we attempt to avoid the 4 dollar parking fee.



Q7)

$$7) a) \quad | \max_a f(a) - \max_a g(a) | \leq \max_a |f(a) - g(a)|$$

case 1: $\max_a f(a) \geq \max_a g(a)$

$$a_f = \arg \max_a f(a)$$

$$a_g = \arg \max_a g(a)$$

$$g(a_f) \leq g(a_g)$$

$$| \max_a f(a) - \max_a g(a) | \leq f(a_f) - g(a_g) \leq f(a_f) - g(a_f)$$

$$| \max_a f(a) - \max_a g(a) | \leq f(a_f) - g(a_f) \quad \text{definition of max}$$

$$\text{for } f(a_f) - g(a_g) = |f(a_f) - g(a_g)| \leq \max_a |f(a) - g(a)|$$

$$\Rightarrow | \max_a f(a) - \max_a g(a) | \leq \max_a |f(a) - g(a)|$$

for case 2: $\max_a g(a) \geq \max_a f(a)$

It would be exactly similar to case 1 with this difference that we need to change the role $f(a)$ and $g(a)$ and we will find

$$| \max_a g(a) - \max_a f(a) | \leq \max_a |g(a) - f(a)|$$

$$7) b) \|BV_i - BV_i'\|_\infty \leq \gamma \|V_i - V_i'\|_\infty$$

$$\|BV_i - BV_i'\|_\infty = \left\| \max_a \sum p(s', r|s, a) [r + \gamma V_i(s')] - \max_a \sum p(s', r|s, a) [r + \gamma V_i'(s')] \right\|_\infty$$

$$= \left\| \max V_i(s') - \max V_i'(s') \right\|_\infty \leq \gamma \max \|V_i(s') - V_i'(s')\|_\infty$$

based on
part a

$$\gamma \max \|V_i(s') - V_i'(s')\|_\infty \leq \gamma \|V_i - V_i'\|_\infty$$

$$\Rightarrow \|BV_i - BV_i'\|_\infty \leq \gamma \|V_i - V_i'\|_\infty$$