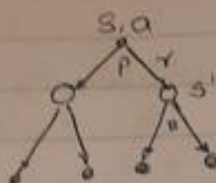


### Ex - 3

i) Action Value Function.



$q_v$  backup diagram.

=)

a) Eq<sup>n</sup> for  $V_\pi$  in  $q_\pi$  &  $\pi$ .

$$V_\pi = \sum_{a \in A(s)} \pi(a|s) q_\pi(s, a)$$

Summation is over all actions  $a$  in the action space  $A(s)$  for the state

$\pi(a|s)$  is the probability of taking action  $a$  in state  $s$  under policy  $\pi$ .

$q_\pi(s, a)$  is expected return starting from state  $s$ , taking action and thereafter following the policy.

b) Eq<sup>n</sup> for  $q_\pi$  in terms of  $V_\pi$  and four-argument  $P$ .

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

$P(s', r | s, a)$  probability of transitioning to  $s'$  receiving reward  $r$  in action  $a$  in  $s$ .

$V_\pi(s')$  value of state  $s'$  under policy  $\pi$ .

c) Bellman eq<sup>n</sup> for action values, for  $q_*$

$$q_*(s, a) = \sum_{s', r} P(s', r | s, a) \left[ r + \gamma \sum_{a'} \pi(a' | s) q_*(s', a') \right]$$

$P(s', r | s, a)$  transition probability function.

2) Fun with Bellman.

=>

a) Eq<sup>n</sup> for  $V_*$  in terms of  $q_*$

$$V_*(s) = \max_a q_*(s, a)$$

b) Eq<sup>n</sup> for  $q_*$  in terms of  $V_*$  four-argument  $P$ .

$$q_*(s, a) = \sum_{s', r} P(s', r | s, a) (r + \gamma V_*(s'))$$

c) Eq<sup>n</sup> for  $\pi_*$  in term of  $q_*$

$$\pi_*(s) = \operatorname{argmax}_a q_*(s, a)$$

d) Eq<sup>n</sup> for  $\pi_*$  in terms of  $V_*$  & four-argument  $P$ .

$$\pi_*(s) = \operatorname{argmax}_a \sum_{s', r} P(s', r | s, a) + (r + \gamma V_*(s'))$$

e) Bellman eq<sup>n</sup> for four value functions ( $V_\pi, V_*, q_\pi, q_*$ ) in terms of three argument function  $P$

$$P(s' | s, a) = \sum_{r \in R} P(s', r | s, a) \quad \text{--- (3.4)}$$

$$r(s, a) = \sum_{s', r} P(s', r | s, a) r \quad \text{--- (3.5)}$$

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

$$\therefore V_*(s) = \max_a [r(s, a) + \sum_{s'} P(s' | s, a) \gamma V_*(s')] ]$$

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

$$\therefore q_*(s, a) = r(s, a) + \sum_{s'} P(s' | s, a) \max_a q_*(s', a)$$

### 3) Fixing Policy Iteration.

- a) While fixing the bug in policy iteration, we can add a termination condition so that the policy remains unchanged for certain no. of iteration. If the policy doesn't change the predetermined no. of iteration, the algorithm terminates.

Changing the judgement condition in pseudo code so that policy to be optimized.

Policy - Stable  $\leftarrow$  true

For each  $s \in S$ :

old action  $\leftarrow \pi(s)$

$\pi(s) \leftarrow \operatorname{argmax}_{s' \rightarrow} p(s', r | s, a)$

$[r + \gamma V(s)]$

If  $q(s', a) \neq q_{\pi}(s, a)$   
old action  
 then policy - stable  $\leftarrow$  false

$\therefore$  The policy then stop & return  $V \approx V_{\pi}$  &  $\pi \approx \pi_{\pi}$   
 else go to 2.

b)  $V_{k+1}(s) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')] \rightarrow$

NO, analogous bug in value iteration as it always selects the best action in each state based on current value fun<sup>n</sup>.

The policy is optimal as long as the result is optimal.



#### 4) Policy iteration for action values:

##### a) 1) Initialization:

$Q(s,a) \in \mathbb{R}$  &  $\pi(s) \in A(s)$  arbitrarily for all  $s \in S$

##### 2) Policy Evaluation:

Loop:  $\Delta \leftarrow 0$

Loop for each  $s \in S$ :

$a \leftarrow \pi(s)$

$q \leftarrow Q(s,a)$

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

$$\Delta \leftarrow \max (\Delta, |q - Q(s,a)|) \quad \left[ \begin{array}{l} \text{until we get} \\ \text{the smallest} \\ \text{positive no.} \end{array} \right]$$

$$\Delta < \theta$$

##### 3) Policy Improvement:

Policy - stable  $\leftarrow$  True

For each  $s \in S$ :

old action  $\leftarrow \pi(s)$

$\pi(s) \leftarrow \operatorname{argmax}_a Q(s,a)$

$$\text{If } \sum_{s',r} p(s',r | s, \pi(s)) \left[ r + \gamma \sum_{a'} \pi(a' | s') Q(a' | s') \right]$$

$\neq$

$$\sum_{s',r} p(s',r | s, \text{old action}) \left[ r + \gamma \sum_{a'} \pi(a' | s') Q(a' | s') \right]$$

$\therefore$  Policy - stable  $\leftarrow$  False

$\Leftarrow$  If policy - stable, then stop & return  $Q \approx q^*$  &  $\pi \approx \pi^*$  else we can do (2)

b) Value iteration update for all action values  
 $q_{k+1}(s,a)$  is  $q_{k+1}(s,a) = \sum_{s',r} p(s',r|s,a) [r + \gamma \max_{a'} q_k(s',a')]$   
 $\theta > 0$

Initialize  
 $Q(s,a)$  for all  $s \in S^*$ , arbitrary except  
 the  $Q$  terminal.

loop:  $\Delta \leftarrow 0$   
 loop for each  $s \in S$   $a \leftarrow \operatorname{argmax}_a Q(s,a)$   
 $q \leftarrow Q(s,a)$

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

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

until  $\Delta < \theta$

output a deterministic policy  $\pi \approx \pi^*$  such that

$$\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma q(s',\pi(s))]$$

$$q_{k+1}(s,a) \leftarrow \sum_{s',r} p(s',r|s,a) [r + \gamma \max_{a'} q_k(s',a')]$$

5)

a), b), c)

```
Answer5a V(s) -
[[ 3.31359559  8.79292942  4.43113177  5.32556099  1.4955287 ]
 [ 1.52582318  2.99591435  2.2534199  1.91064941  0.55045095]
 [ 0.05486787  0.74165922  0.67626363  0.36114423 -0.40025498]
 [-0.96965064 -0.43208514 -0.35180898 -0.58272448 -1.18027658]
 [-1.85380443 -1.34185832 -1.22622928 -1.42007309 -1.97241846]]

Answer5b V(s) -
[[21.9764967  24.41877948 21.97690153 19.41877948 17.47690153]
 [19.77884703 21.97690153 19.77921138 17.80129024 16.02116122]
 [17.80096232 19.77921138 17.80129024 16.02116122 14.4190451 ]
 [16.02086609 17.80129024 16.02116122 14.4190451  12.97714059]
 [14.41877948 16.02116122 14.4190451  12.97714059 11.67942653]]

Pi(s) -
[['right' 'up' 'left' 'up' 'left']
 ['up' 'up' 'up' 'left' 'left']
 ['up' 'up' 'up' 'up' 'up']
 ['up' 'up' 'up' 'up' 'up']
 ['up' 'up' 'up' 'up' 'up']]

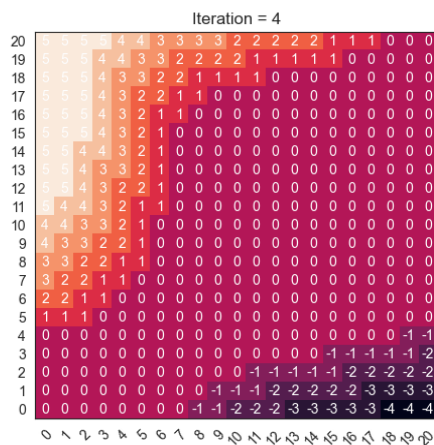
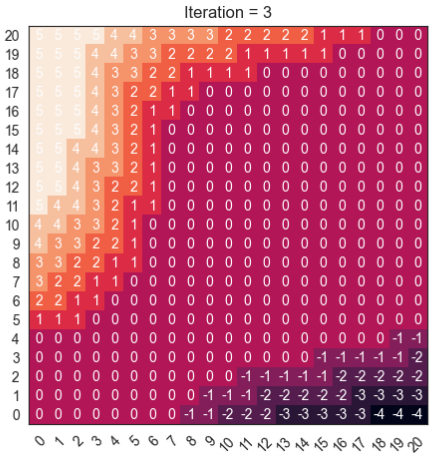
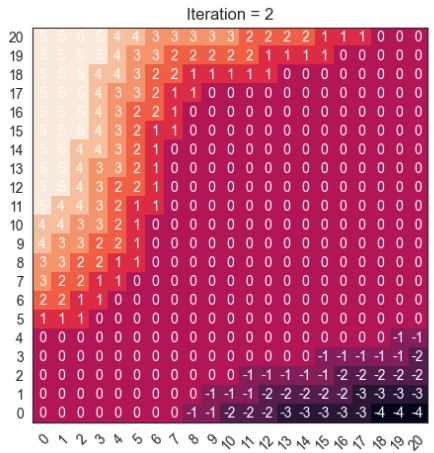
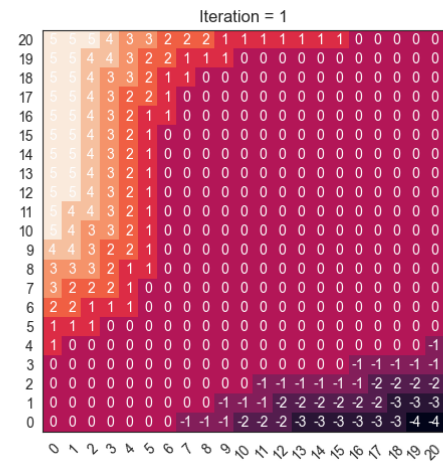
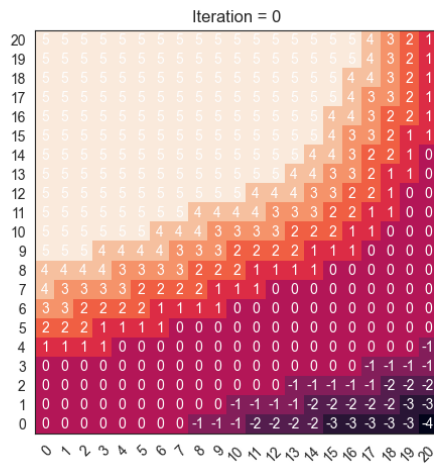
Answer5c V(s) -
[[21.97748529 24.4194281  21.97748529 19.4194281  17.47748529]
 [19.77973676 21.97748529 19.77973676 17.80176308 16.02158677]
 [17.80176308 19.77973676 17.80176308 16.02158677 14.4194281 ]
 [16.02158677 17.80176308 16.02158677 14.4194281  12.97748529]
 [14.4194281  16.02158677 14.4194281  12.97748529 11.67973676]]

Pi(s) -
[['right' 'up' 'left' 'up' 'left']
 ['up' 'up' 'up' 'left' 'left']
 ['up' 'up' 'up' 'up' 'up']
 ['up' 'up' 'up' 'up' 'up']
 ['up' 'up' 'up' 'up' 'up']]

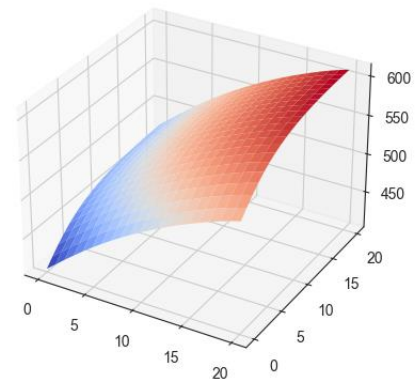
□
```

6)

- a. In these four graphs, an off-white color represents positive values, while purple denotes negative values. The color at the midpoint of each graph's color scale signifies a value of '0'.



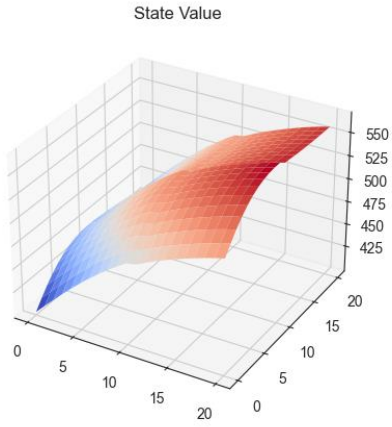
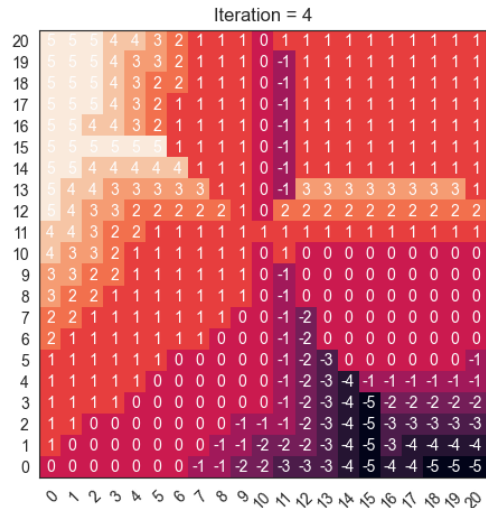
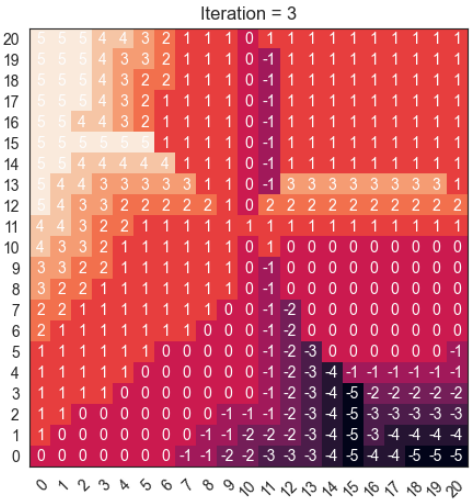
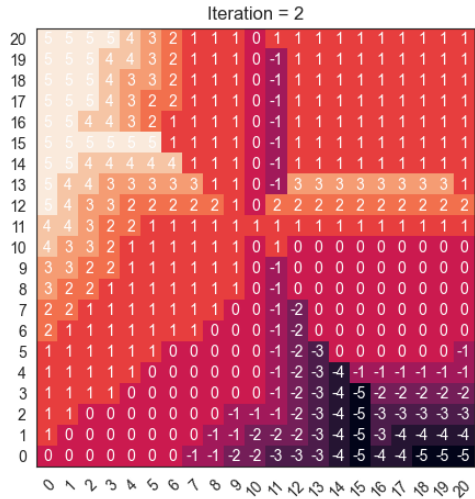
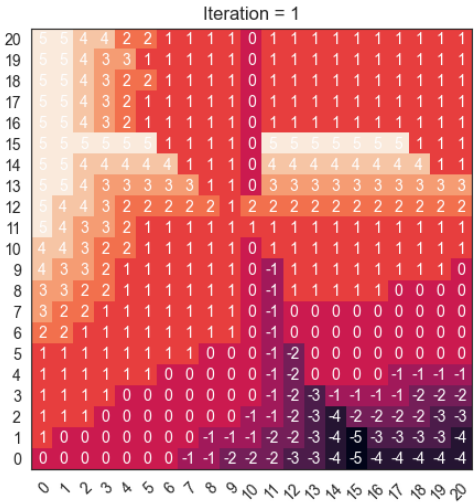
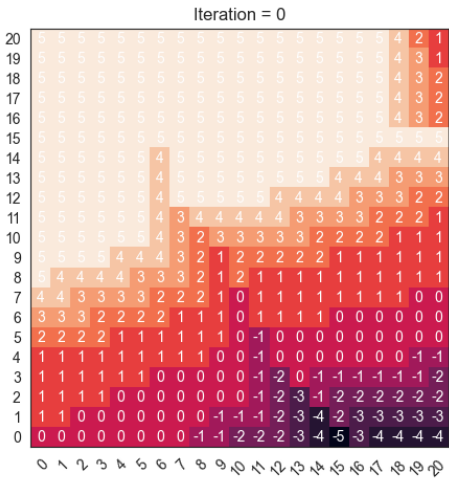
State Value



- b. When modifications are implemented to the car rental system, they influence the reward structure. Specifically, if an employee helps in relocating cars, it decreases



the cost associated with moving cars in one direction by a unit. This implies that the employee lives near the second location, resulting in a modified reward for transferring cars from the first location to the second.



The adjustments to the car rental system's reward structure incentivize strategic car movements and inventory management: Employee Assistance: Moving a car from lot 1 to lot 2 is cheaper by one unit, encouraging movement towards lot 2 to save costs. Parking Fee: An added parking fee for more than 10 cars at any location discourages excessive inventory, promoting a strategy to keep car counts close to 10 at each location to maximize rentals while avoiding fees. Overall, these changes favor strategies that balance car availability with minimizing costs, shifting optimal operations towards maintaining moderate inventory levels at both locations to optimize profitability.

7) Extra Credit:

a) For any fun<sup>n</sup>  $f$  and  $g$ .

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

if  $\max_a f(a) - \max_a g(a) \geq 0$

$$\left| \max_a f(a) - \max_a g(a) \right| = \max_a f(a) - \max_a g(a)$$

$\max_a g(a) \geq g(x)$  for all  $x$ .

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

for all  $x$ .

$$\max_a f(a) - \max_a g(a) \leq \max_a (f(a) - g(a)) \quad \text{--- ①}$$

if,

$a_1 = \arg\max f(a)$ ,  $a_2 = \arg\max [f(a) - g(a)]$   
then

$$f(a_1) - g(a_1) \leq f(a_2) - g(a_2) \quad \text{--- ②}$$

$$f(a_1) - g(a_1) = \max_a (f(a) - g(a)) \quad \left. \vphantom{\max_a (f(a) - g(a))} \right\} \text{② in ①}$$

$$f(a_2) - g(a_2) = \max_a (f(a) - g(a))$$

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

When,  $\max_a f(a) - \max_a g(a) < 0$

Substituting  $f(a) \leftarrow g(a)$  &  $g(a) \rightarrow f(a)$ .

$$\therefore \left| \max_a f(a) - \max_a g(a) \right| \leq \max_a |f(a) - g(a)| \quad \text{--- (3)}$$

b) Let  $B$  denote the Bellman backup operator,  
Value iteration

$\hookrightarrow$  Vector form :  $V_{K+1} \leftarrow BV_K$

$$\begin{aligned} |\beta V_i(s) - \beta V_i(s')| &= \left| \max_a \sum_{s',r} p(s',r|s,a) [r + \beta V_i(s')] \right. \\ &\quad \left. - \max_a \sum_{s',r} p(s',r|s',a) [r + \beta V_i(s')] \right| \end{aligned} \quad \text{--- (4)}$$

Taking (3) & from a) we get,

$$\begin{aligned} \text{eqn 4} &\leq \max_a \left| \sum_{s',r} p(s',r|s,a) [r + \gamma V_K(s')] \right. \\ &\quad \left. - \sum_{s',r} p(s',r|s',a) [r + \gamma V_K(s')] \right| \\ &= \max_a \left| \sum_{s',r} p(s',r|s,a) \times \gamma (V_K(s') - V_K(s)) \right| \\ &\quad \text{for all } s \in S \end{aligned}$$



for  $n$ -length  $v_i$  and  $v_i'$ .

$$\| \beta v_i - \beta v_i' \| = \max \left\{ \begin{array}{l} | \beta v_i(s_1) - \beta v_i'(s_1) | \\ | \beta v_i(s_2) - \beta v_i'(s_2) | \\ \vdots \\ | \beta v_i(s_n) - \beta v_i'(s_n) | \end{array} \right\} \quad \text{--- (5)}$$

$$\text{eqn (5)} \leq \left\{ \begin{array}{l} \max_a \left| \sum_{s,r} p(s,r|s,a) \gamma \| v_i - v_i' \| \right| \\ \max_a \sum_{s,r} p(s,r|s,a) \gamma \| v_i - v_i' \| \end{array} \right\}$$

$$= \max \left\{ \begin{array}{l} \gamma \| v_i - v_i' \| \\ \gamma \| v_i - v_i' \| \\ \vdots \\ \gamma \| v_i - v_i' \| \end{array} \right\} = \boxed{\gamma \| v_k - v_k' \|}$$

c) Banach Fixed point theorem  $\Rightarrow d(T(x), T(y)) \leq q d(x, y)$ .

$$\Rightarrow d(x, y) = \| x - y \|$$

$T(x) = \beta x$ ;  $\beta$  is contracting mapping  
 $q = \gamma$ .

$$\text{from (b)} \| \beta v_i - \beta v_i' \|_{\infty} \leq \gamma \| v_i - v_i' \|_{\infty} \quad \text{--- (6)}$$

$$n = \| v_{n+1} - v_n \| \leq \gamma^n \| v_1 - v_n \|$$

for any  $m, n$ .

$$\| v_m - v_n \| \leq \| v_m - v_{m-1} \| + \| v_{m-1} - v_{m-2} \| + \dots + \| v_{n+1} - v_n \|$$

$$\leq \gamma^{m-1} \|v_1 - v_0\| + \gamma^{m-2} \|v_1 - v_0\| \dots$$

$$\gamma^m \|v_1 - v_0\|$$

$$\leq \gamma^m \|v_1 - v_0\| \sum_{k=0}^{\infty} \gamma^k$$

$$\rightarrow \|v_m - v_0\| \leq \gamma^m \|v_1 - v_0\| \frac{1}{1-\gamma} < \varepsilon$$

$\varepsilon > 0$  is an arbitrary value for large  $N$   
here  $\text{seq}^n$ .

$v_i$  is Cauchy, matching it a fixed point.

$$V^* = \lim_{n \rightarrow \infty} v_n = \lim_{n \rightarrow \infty} \beta(v_{n-1}) = \beta(\lim_{n \rightarrow \infty} v_{n-1})$$

$$= \beta V^*$$

Coverage to fix point -  
 $v_1, v_2$  then.

$$v_1 = \beta v_1, v_2 = \beta v_2$$

$$\|\beta v_1 - \beta v_2\|_{\infty} = \|v_1 - v_2\|_{\infty} \leq \gamma \|v_1 - v_2\|_{\infty}$$

$$\forall \varepsilon \in (0, 1)$$

$$\|v_1 - v_2\| = 0 \Rightarrow \text{unique fix point.}$$

Eq. 3.9

$$\beta v_A(s) = \max_a E[R_{t+1} + \gamma v_A(s_{t+1}) | s_t = s, a_t = a]$$

$$\Rightarrow \beta v_A(s) = \max_a q_{\pi^*}(s, a)$$

$$= v_{\pi^*}(s)$$

$\therefore$  The this unique fixed point is equiv to Bellman.