

* Exercise 3: Dynamic Programming *

Discussed with: Vaibhav Garaf

Q.1) a) Give an equation of v^* in terms of q^*

$$\rightarrow v^*(s) = \sum_a \pi^*(a|s) q^*(s, a)$$

b) equation of q^* in terms of v^* & four argument P

$$\rightarrow q^*(s, a) = \sum_{s', r} P(s', r|s, a) (r + \gamma v^*(s'))$$

c) equation of π^* in terms of q^*

$$\rightarrow \pi^*(s) = \arg\max_a q^*(s, a)$$

d) equation of π^* in terms of v^* & four argument P

$$\rightarrow \pi^*(s) = \arg\max_a \sum_{s', r} P(s', r|s, a) [r + \gamma v^*(s')]$$

e) rewrite 4 bellman equations for (v_π, v^*, q_π, q^*) in terms of 2 args function P & 2 args function r

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

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

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

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

Q.2) Fixing policy iteration.

a)

→ we can fix the bug in policy iteration by every time selecting optimal actions with equal probability.

for example:- If a policy consists of 4 actions

$$a = \{ \text{left, Right, Up, down} \}$$

at each iteration we will choose optimal policy with equal probability.
 if left & up are optimal,
 $\therefore \pi(a) = \{0.5, 0, 0.5, 0\}$

∴ Modified pseudocode:-

1. Initialization:-

$$V(s) \in \mathbb{R} \quad \& \quad \pi(s) \in A(s) \quad \forall s \in S$$

2. Policy Evaluation:-

loop :-

$$\Delta \leftarrow 0$$

loop for each $s \in S$:

$$U \leftarrow V(s)$$

$$V(s) \leftarrow \sum_{s', r} p(s', r | s, \pi(s)) [\alpha + U(s')]$$

$$\Delta \leftarrow \max(\Delta, |U - V(s)|)$$

until $\Delta < \epsilon$ (ϵ a small positive number)

3. Policy improvement

policy-stable \leftarrow True

for each $s \in S$:

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

$$\pi(s) \leftarrow \underset{\text{no-optimal actions}}{1} * \arg\max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$$

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

If policy-stable, then stop & return V^* & $\pi^* \approx \pi$; else go to 2.

b) I don't think that value iteration has any bug similar to policy iteration.

When we find optimal value function it should be unique, regardless of how many different optimal policies we have.

Q.3) Policy iteration for action values

a) Initialization:

$Q(s, a) \in R$, $\pi(s) \in A(s)$ arbitrarily
 $\forall s \in S$

Policy evaluation

loop:-

$$\Delta \leftarrow 0$$

loop for each $s \in S$:-

loop for each $a \in A$:-

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

$$\Delta \leftarrow \max_a (\Delta, \text{abs}(q(s, a) - Q(s, a)))$$

until $\Delta < 0$

Policy improvement :-

policy_stable \leftarrow true

for each $s \in S$:-

for each $a \in A$:-

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

if oldaction $\neq \pi(s)$

policy stable \leftarrow false

If policy stable then stop & return

$$Q^{\pi^*} \text{ & } \pi^*$$

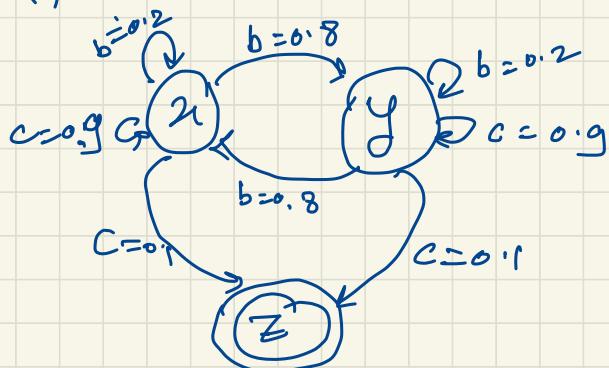
else go to policy evaluation

b)

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

Q. h) Policy iteration by hand

a)



→ agent's goal is to reach state z (terminal state)

→ But from any state x & y probability of going to terminal state z by taking action c' is equal to 0.1.

→ Thus now we will need to look at the rewards

→ On state x agent receives reward '-1' & on state y '-2'

→ To minimize negative rewards agent should stay at x

→ This optimal policy could be everytime, from y → move to x → take c from x from x → take c

b) Iteration 0

$$\pi_0 = [x \rightarrow c, y \rightarrow c]$$

$$V(x) = r + \gamma \sum_{s'} p(s'|x, c) V(s')$$
$$= -1 + 0.1 V(z) + 0.9 V(y)$$

$$V(y) = -2 + 0.1 V(z) + 0.9 V(y)$$

$$V(z) = 0$$

$$0.1 V(x) = -1 \quad = -10$$
$$0.1 V(y) = -2 \quad = -20$$

Policy improvement

$$\pi(x) = \arg\max_a \left[r + \gamma \sum_s p(s|x, c) V(s) \right]$$

$$= \arg\max_a \begin{cases} -1 + 0.8 * (-20) + 0.2 * (-10), \\ -1 + 0.9 * (-10) \end{cases}$$

$$= \arg\max_a \begin{cases} -1 - 16 - 2 = -19, & \leftarrow b \\ -1 - 9 = -10 & \leftarrow c \end{cases}$$

$$\pi(x) = c$$

$$\pi(y) = \arg\max_a \begin{cases} -2 + 0.8 * (-10) + 0.2 * (-20), \\ -2 + 0.1 * 0 + 0.9 * (-20) \end{cases}$$

$$= \arg\max_a \{-14, -20\} = b$$

$$\therefore \Pi_1 = [x \rightarrow c, y \rightarrow b]$$

Iteration 1

$$V(x) = -1 + 0.1 * V(z) + 0.9 * V(x)$$

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

$$V(z) = 0$$

$$\therefore V(x) = -10 \quad V(y) = \frac{-10}{0.8} = -12.5$$

Policy improvement.

$$\pi(x) = \arg \max_a \left\{ \begin{array}{l} -1 + 0.8 * (-12.5) + 0.2 * (-10) \\ -1 + (0.1 * 0) + 0.9 * (-10) \end{array} \right\},$$

$$= \arg \max_a \left\{ \begin{array}{l} -1 - 10 - 2 \\ -1 - 9 \end{array} \right\},$$

$$= \arg \max_a \left\{ \begin{array}{l} -13 \xleftarrow{b} \\ -10 \xleftarrow{c} \end{array} \right\}$$

$$\pi(x) = c$$

$$\pi(y) = \arg \max_a \left\{ \begin{array}{l} -2 + 0.8 * (-10) + 0.2 * (-12.5) \\ -2 + 0.9 * (-12.5) \end{array} \right\},$$

$$\pi(y) = \arg \max_a \left\{ \begin{array}{l} -12.5 \xleftarrow{b} \\ -13.25 \xleftarrow{c} \end{array} \right\}$$

$$\pi(y) = b$$

$$\therefore \pi = [x \rightarrow c, y \rightarrow b]$$

which is same as previous policy. Hence policy is converged.

C)

Initial policy

$$\pi = [x \rightarrow b, y \rightarrow b]$$

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

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

$$v(z) = 0$$

$$0.8 v(x) = -1 + 0.8 v(y)$$

$$0.8 v(y) = -2 + 0.8 v(x)$$

$$0.8 v(x) = -1 + [-2 + 0.8 v(x)]$$

$$0.8 v(x) = -3 + 0.8 v(x)$$

$$0 = 3 \quad \text{which is inconsistent.}$$

Our initial policy always takes action b which causes it to loop between state x & y infinitely.

Since we are considering undiscounted MDP, we will get sum of infinite negative numbers.

Does discounting help?

- Discounting factor can help with above said problem as values in future will become insignificant
- But ensuring accurate optimal solution will be difficult in this case as the overall problem changes.
- In this MDP Optimal policy is not dependent on discounting factor, as the policy will still be able to converge even with discount rate.

Q.7) Proving Convergence of value iteration

for any function $f \& g$

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

→ case 1

$$\max_a f(a) - \max_a g(a) > 0 \quad -\textcircled{1}$$

$$\text{let } \max_a f(a) = a_f^*$$

$$\max_a g(a) = a_g^*$$

$$\max_a f(a) - \max_a g(a) = f(a_f^*) - g(a_g^*) > 0$$

-\textcircled{2}

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

From \textcircled{2} we can say that

$g(a_f^*) \leq g(a_g^*)$ as a_g^* is optimal action for g .

Thus we can write inequality from ① as.

$$\begin{aligned} f(a_f^*) - g(a_f^*) &> f(a_f^*) - g(a_g^*) \geq 0 \\ \therefore f(a_f^*) - g(a_g^*) &= |f(a_f^*) - g(a_f^*)| \\ \therefore f(a_f^*) - g(a_g^*) &\leq f(a_f^*) - g(a_f^*) \\ &= |f(a_f^*) - g(a_f^*)| \leq \max_a |f(a) - g(a)| \end{aligned}$$

Case 2:-

$$\max_a f(a) - \max_a g(a) = f(a_f^*) - g(a_g^*) < 0$$

$$f(a_g^*) \leq f(a_f^*)$$

Considering above

$$g(a_g^*) - f(a_g^*) > g(a_g^*) - f(a_f^*) > 0$$

as function

$$g(a_g^*) - f(a_g^*) = |g(a_g^*) - f(a_g^*)| \text{ } g > \text{func-f}$$

$$\therefore g(a_g^*) - f(a_f^*) \leq g(a_g^*) - f(a_g^*)$$

$$= |g(a_g^*) - f(a_g^*)| \leq \max_a |g(a) - f(a)|$$

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

Q.7) b)

From Q7 part a)

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

$U_k \leftarrow$ k-th value iteration

$\mathcal{B} \leftarrow$ bellman backup operator.

$$U_{k+1} \leftarrow \mathcal{B}U_k$$

$$U_{k+1}(s) \leftarrow \max_a \sum p(s'|s, a) [r + \gamma U_k(s')]$$

$\|U\|_\infty$ is max absolute value of its elements



From definition

$$\begin{aligned} \| \mathcal{B}U_k - \mathcal{B}U_k' \|_\infty &= \| R(s) + \max_a \gamma \sum p(s'|s, a) V_{k+1}(s') \\ &\quad - (R(s) + \max_a \gamma \sum p(s'|s, a) V'(s')) \|_\infty \end{aligned}$$

$$= \max_a \gamma \sum p(s'|s, a) \| V(s') - \max_a \gamma \sum p(s'|s, a) V'(s') \|_\infty$$

= From previous, above eqn. is less than or equal to,
eqn.

$$\leq \max_a \left| \gamma \sum p(s'|s, a) V(s') - \gamma \sum p(s'|s, a) V'(s') \right|$$

$$\leq \max_a \gamma \left| \sum p(s'|s, a) V(s') - \sum p(s'|s, a) V'(s') \right|$$

$$\leq \max_a \gamma \sum s' P(s'|s,a) \|v - v''\|$$

$$= \max_a \gamma \|v - v'\|$$

We know that sum of probabilities is given.
A state & action to take

$$\sum_{s'} P(s'|s,a) = 1$$

$$\therefore \max_a \gamma \|v - v'\| = \gamma \|v - v'\|$$

$$\therefore \|Bv_i - Bv_i''\|_\infty \leq \gamma \|v_i - v_i''\|_\infty$$

Question 5 Answer

Part A

Value Function for iterative policy evaluation

$$[-1.9 \quad -1.3 \quad -1.2 \quad -1.4 \quad -2. \quad -1. \quad -0.4 \quad -0.3 \quad -0.6 \quad -1.2 \quad 0.1 \quad 0.7 \quad 0.7 \quad 0.4 \\ -0.4 \quad 1.5 \quad 3. \quad 2.3 \quad 1.9 \quad 0.6 \quad 3.3 \quad 8.8 \quad 4.4 \quad 5.3 \quad 1.5]$$

Part B

Optimal Value function for Value Iteration

```
[14.4 16. 14.4 13. 11.7 16. 17.8 16. 14.4 13. 17.8 19.8 17.8 16.  
14.4 19.8 22. 19.8 17.8 16. 22. 24.4 22. 19.4 17.5]
```

```
("Optimal policy for Value Iteration default dict(<class 'list'>, {(0, 0): "
  "[ 'RIGHT', 'UP' ], (0, 1): [ 'UP' ], (0, 2): [ 'LEFT', 'UP' ], (0, 3): [ 'LEFT', "
  "'UP' ], (0, 4): [ 'LEFT', 'UP' ], (1, 0): [ 'RIGHT', 'UP' ], (1, 1): [ 'UP' ], (1, "
  "2): [ 'LEFT', 'UP' ], (1, 3): [ 'LEFT', 'UP' ], (1, 4): [ 'LEFT', 'UP' ], (2, 0): "
  "[ 'RIGHT', 'UP' ], (2, 1): [ 'UP' ], (2, 2): [ 'LEFT', 'UP' ], (2, 3): [ 'LEFT', "
  "'UP' ], (2, 4): [ 'LEFT', 'UP' ], (3, 0): [ 'RIGHT', 'UP' ], (3, 1): [ 'UP' ], (3, "
  "2): [ 'LEFT', 'UP' ], (3, 3): [ 'LEFT' ], (3, 4): [ 'LEFT' ], (4, 0): [ 'RIGHT' ], "
  "(4, 1): [ 'LEFT', 'DOWN', 'RIGHT', 'UP' ], (4, 2): [ 'LEFT' ], (4, 3): [ 'LEFT', "
  "'DOWN', 'RIGHT', 'UP' ], (4, 4): [ 'LEFT' ]})")
```

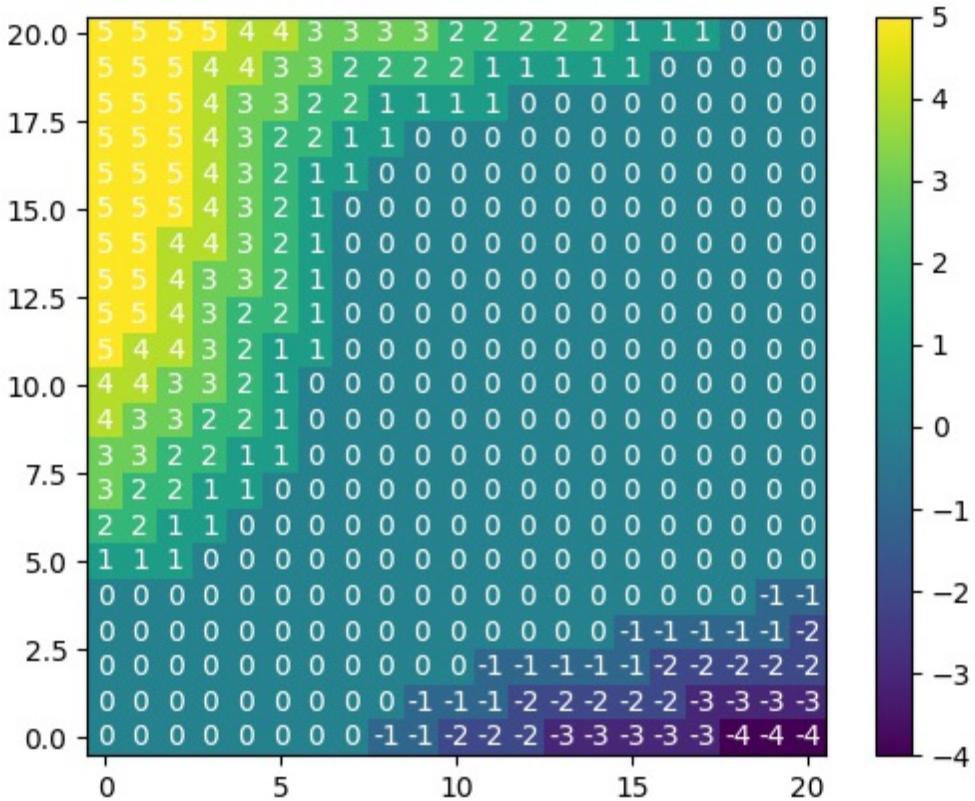
Part C sequence of actions in array is [LEFT,DOWN,RIGHT,UP]

```

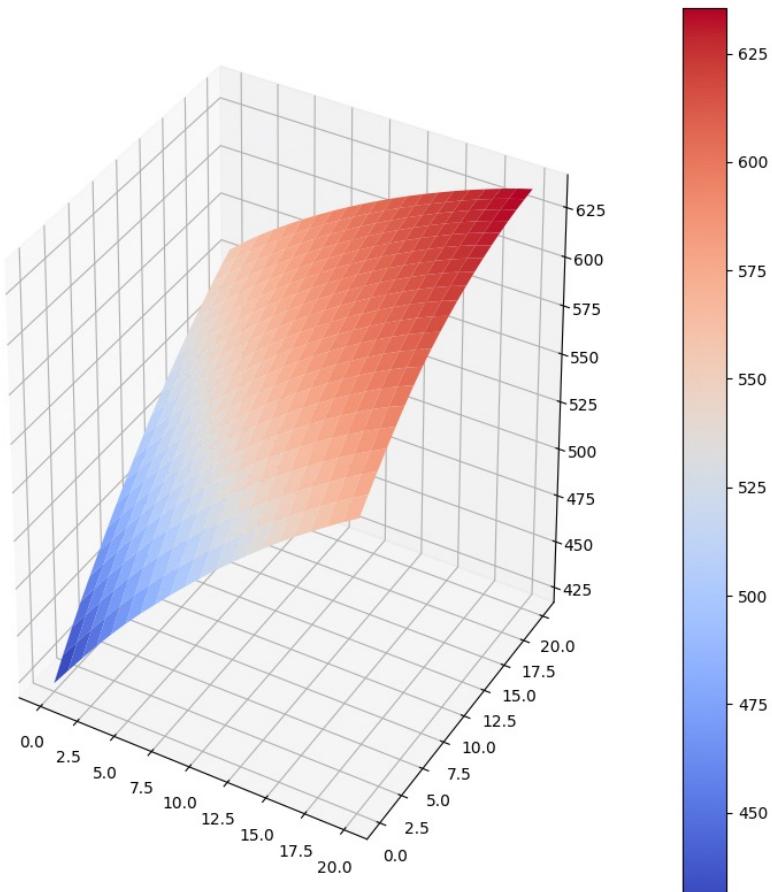
Optimal Value function for Policy Iteration
[14.4 16. 14.4 13. 11.7 16. 17.8 16. 14.4 13. 17.8 19.8 17.8 16.
14.4 19.8 22. 19.8 17.8 16. 22. 24.4 22. 19.4 17.5]
('Optimal policy function for Policy Iteration \n'
[ [0. 0. 0.5 0.5 ]\n'
[ [0. 0. 0. 1. ]\n'
[ [0.5 0. 0. 0.5 ]\n'
[ [0. 0. 0.5 0.5 ]\n'
[ [0. 0. 0. 1. ]\n'
[ [0.5 0. 0. 0.5 ]\n'
[ [0. 0. 0. 0.5 ]\n'
[ [0. 0. 0. 1. ]\n'
[ [0.5 0. 0. 0.5 ]\n'
[ [0.5 0. 0. 0.5 ]\n'
[ [0.5 0. 0. 0.5 ]\n'
[ [0. 0. 0.5 0.5 ]\n'
[ [0. 0. 0. 1. ]\n'
[ [0.5 0. 0. 0.5 ]\n'
[ [0.5 0. 0. 0.5 ]\n'
[ [0.5 0. 0. 0.5 ]\n'
[ [0. 0. 0. 0.5 ]\n'
[ [0. 0. 0. 1. ]\n'
[ [0.5 0. 0. 0.5 ]\n'
[ [1. 0. 0. 0. ]\n'
[ [1. 0. 0. 0. ]\n'
[ [0. 0. 1. 0. ]\n'
[ [0.25 0.25 0.25 0.25 ]\n'
[ [1. 0. 0. 0. ]\n'
[ [0.25 0.25 0.25 0.25 ]\n'
[ [1. 0. 0. 0. ]\n'

```

Question 6 Part A Optimal Policy

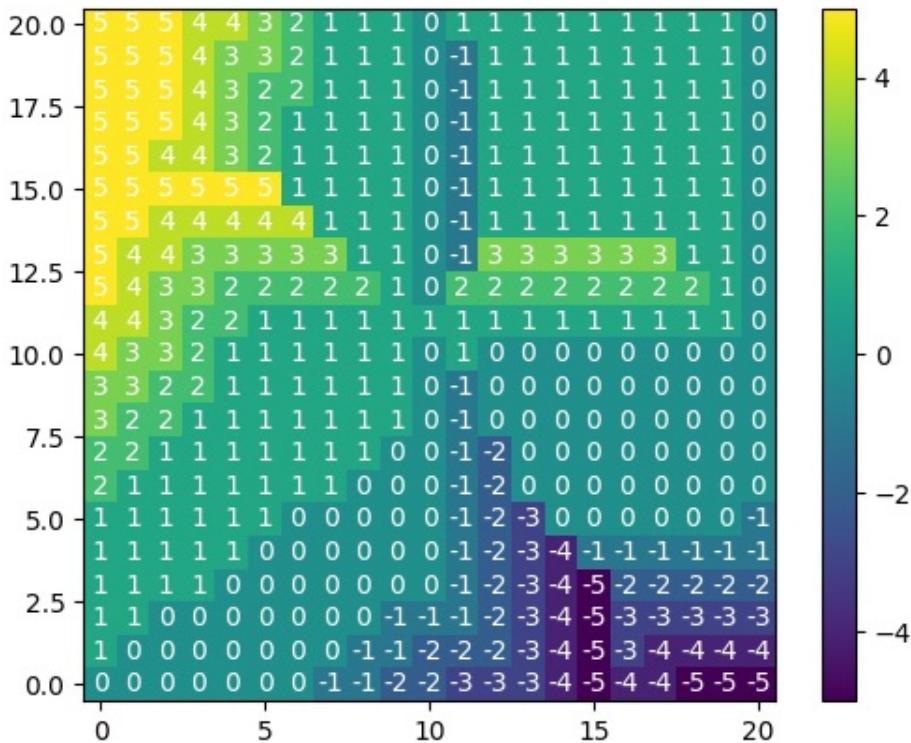


Question 6 A Optimal Value function

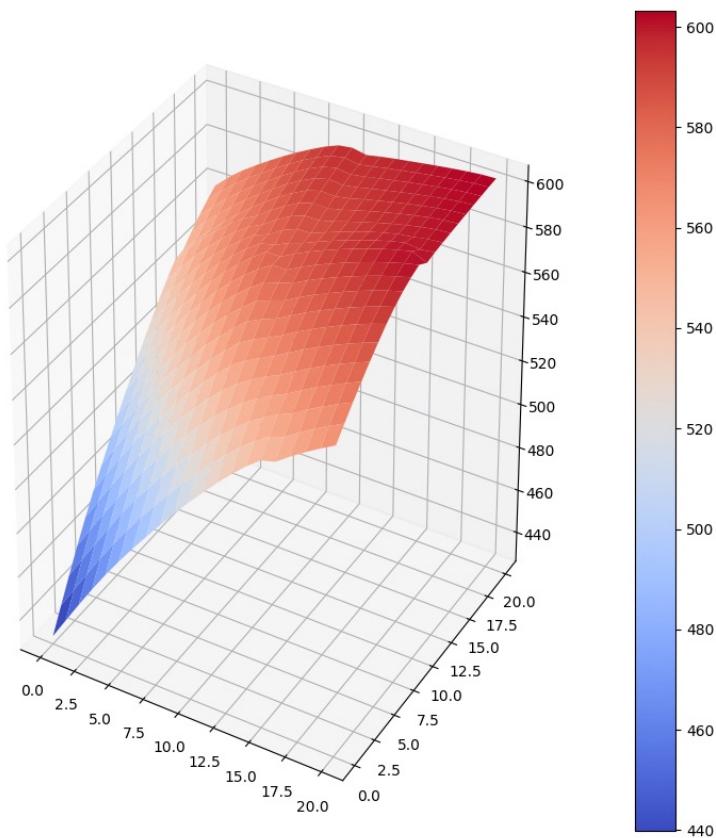


Question 6 Part B (Modified Jacks Car Rental)

Optimal Policy



Optimal Value Function for Modified Jack's Car Rental



written: How do we change dynamics function to reflect modified Jack's car rental problem

Answer:

We mainly change the way we calculate reward to ensure we include cases with changed costs. These changes are mainly in overnight cost to move cars. We check if the action is positive, i.e. moving car from A to B, then our cost is reduced by 2\$ since we get 1 free car move. Also, we check if at any location we have more than 10 cars. Thus, we place additional cost of 4\$.

Written: How does final policy differ from original Jack's Car Rental Problem

Ans :

One immediate observation can be done in optimal policy of jack's car rental is, it tries to move atleast 1 car from location 1 to 2 since its free to move except when location 2 has exactly 10 cars, since we'll need to pay overflow parking cost