

Submitted by: Aadhar Bawali

Discussed with: Christopher Anton Dominic
Tharvan Sanghvi

1. 1 point. (RL2e 3.25 - 3.29) Fun with Bellman.

Written:

- (a) Give an equation for v_* in terms of q_* .
 (b) Give an equation for q_* in terms of v_* and the four-argument p .
 (c) Give an equation for π_* in terms of q_* .
 (d) Give an equation for π_* in terms of v_* and the four-argument p .
 (e) Rewrite the four Bellman equations for the four value functions (v_π, v_*, q_π, q_*) in terms of the three-argument function p (Equation 3.4) and the two-argument function r (Equation 3.5).

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

b) $q_*(s, a) = \sum_{s' \in S} p(s', r|s, a) [r + \gamma V_*(s')]$

c) $\pi_*(s) = \arg \max_a q_*(s, a)$

d) $\pi_* = \arg \max_a \sum_{s' \in S} p(s', r|s, a) [r + \gamma V_*(s')]$

e) Bellmann eqⁿ

$$V_*(s) = \sum_a \pi(a|s) \sum_{s' \in S} p(s', r|s, a) [r + \gamma V_*(s')] \quad \text{--- (3.4)}$$

$$p(s'|s, a) \doteq \Pr\{S_t=s' | S_{t-1}=s, A_{t-1}=a\} = \sum_{r \in R} p(s', r|s, a). \quad \text{--- (3.5)}$$

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

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

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

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

Q. 2

2. 1 point. (RL2e 4.5, 4.10) Policy iteration for action values.

Written:

- (a) How would policy iteration be defined for action values? Give a complete algorithm for computing q_* , analogous to that on page 80 for computing v_* . Please pay special attention to this exercise, because the ideas involved will be used throughout the rest of the book.
 (b) What is the analog of the value iteration update Equation 4.10 for action values, $q_{k+1}(s, a)$?

Initialization:

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

Policy Evaluation:

Loop:

Δ ← 0

Loop for each $s \in S$ & $a \in A$

$$q_s(a) = \delta(s, a)$$

$$q_s(a) \leftarrow \sum_{s' \in S} p(s'|s, a) \left[r + \gamma \sum_a \pi(a|s) q_{s'}(a) \right]$$

$$\Delta \leftarrow \max(\Delta, |q_s(a) - q_{s'}(a)|)$$

until $\Delta < \delta$ (a small positive number determining the accuracy of estimation)

Policy Improvement

Policy-stable ← true

For each $s \in S$ & $a \in A$:old-action ← $\pi(s)$ $\pi(s) \leftarrow \arg \max_a q_s(a)$ if old-action ≠ $\pi(s)$, then policy-stable ← falseif policy-stable, then stop return $\delta = q_*$ & $\pi = \pi_*$ else go to Policy Evaluation

2(b)

$$q_{k+1}(s, a) = E \left[R_{t+1} + \max_{a'} \gamma q_k(s', a') \right]$$

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

Q. 3

3. 2 points. Policy iteration by hand.

Written: Consider an undiscounted MDP having three states, x, y, z . State z is a terminal state. In states x and y there are two possible actions, b and c . The transition model is as follows:

- In state x , action b moves the agent to state y with probability 0.8 and makes the agent stay put (at state x) with probability 0.2.

- In state y , action b moves the agent to state x with probability 0.8 and makes the agent stay put (at state y) with probability 0.2.

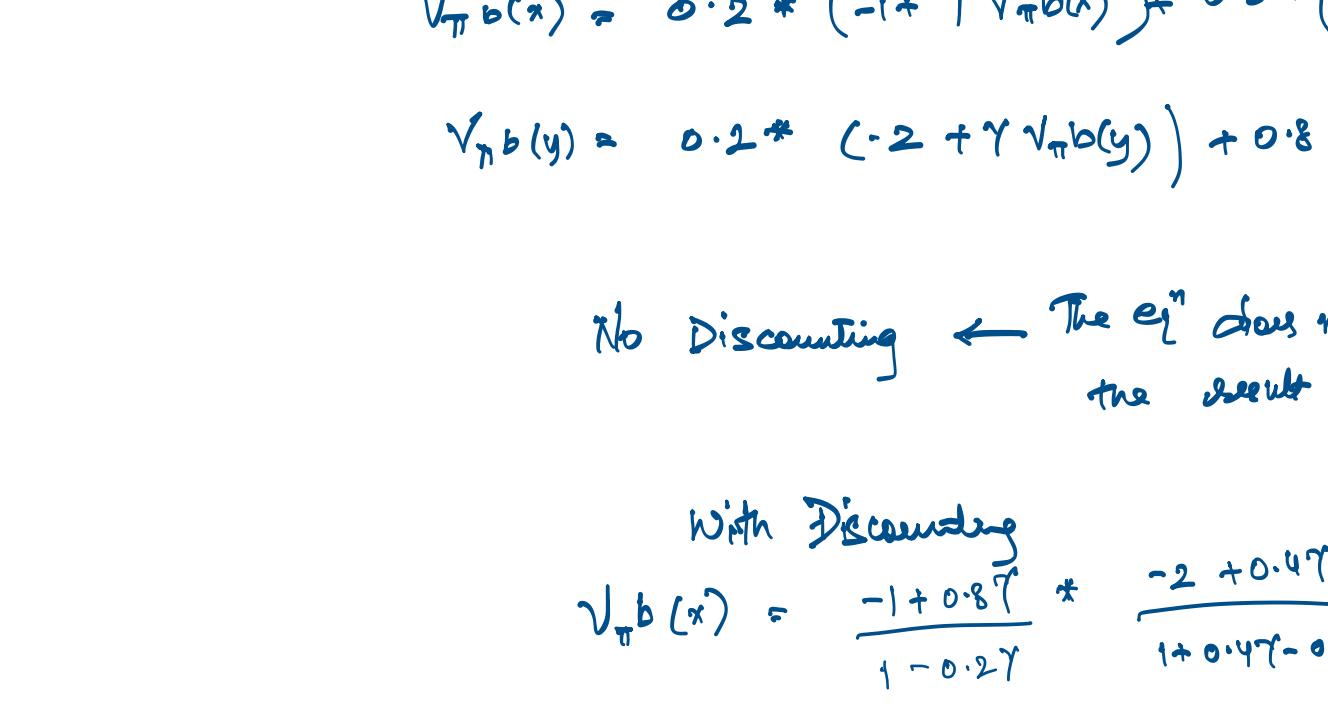
- In either state x or state y , action c moves the agent to state z with probability 0.1 and makes the agent stay put with probability 0.9.

1

The reward model is as follows:

- In state x , the agent receives reward -1 regardless of what action is taken and what the next state is.

- In state y , the agent receives reward -2 regardless of what action is taken and what the next state is.

 $\{x, y, z\}$
 $z = \text{terminal state}$ 

Answer the following questions:

- (a) What can be determined qualitatively about the optimal policy in states
- x
- and
- y
- (i.e., just by looking at the transition and reward structure, without running value/policy iteration to solve the MDP)?

- (b) Apply policy iteration, showing each step in full, to determine the initial policy and the values of states
- x
- and
- y
- . Assume that the initial policy has action
- c
- in both states.

- (c) What happens to policy iteration if the initial policy has action
- b
- in both states? Does discounting help? Does the optimal policy depend on the discount factor (in this particular MDP)?

- a) As reward for terminal state
- z
- is 0
-
- it's better if
- x
- &
- y
- tries to achieve it

for y moving to state x is also beneficial
as the reward for x is -1 & for y is -2 As the action $(y) \xrightarrow{b} (x)$ (less -ve reward state)
∴ it is better for y to have action b & for x taking action would be problematic& c would be better option∴ policy would be choosing action b for y & action c for x

b) Policy evaluation

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

Assumption: Action c in both state

$$V_{\pi_c}(x) = 0.1 * (-1 + 0) + 0.9 * (-1 + V_{\pi_c}(x))$$

$$V_{\pi_c}(y) = 0.1 * (-2 + 0) + 0.9 * (-2 + V_{\pi_c}(y))$$

$$V_{\pi_c}(z) = -0.1 + (-0.9 + 0.9 V_{\pi_c}(z))$$

$$V_{\pi_c}(z) = -0.1 + (-1.8 + 0.9 V_{\pi_c}(z))$$

$$\downarrow -1 = 0.1 V_{\pi_c}(x)$$

$$V_{\pi_c}(x) = -10 \quad \pi(x) \rightarrow c$$

$$-2 = 0.1 V_{\pi_c}(y) \quad \pi(y) \rightarrow c$$

$$V_{\pi_c}(y) = -20$$

Policy improvement

when applying policy π

$$V_{\pi_b}(x) = 0.2 * (-1 + 0) + 0.8 * (-1 + V_{\pi_b}(x))$$

$$V_{\pi_b}(y) = 0.2 * (-2 + 0) + 0.8 * (-2 + V_{\pi_b}(y))$$

$$V_{\pi_b}(z) = -0.1 + (-1.8 + 0.9 V_{\pi_b}(z)) = -1.9$$

$$V_{\pi_b}(z) = -1.8 + 0.9 V_{\pi_b}(z) = -1.8$$

$$\therefore V_{\pi_b}(x) = -1.9 < V_{\pi_c}(x) = -10$$

$$V_{\pi_b}(y) = -1.8 > V_{\pi_c}(y) = -20$$

∴ $\pi'(x) \rightarrow b$ & $\pi'(y) \rightarrow b$

$$\pi'(y) \rightarrow b$$

 $\pi' \neq \pi \rightarrow$ again Policy Evaluation

Policy Evaluation

$$V_{\pi_b}(x) = 0.1 * (-1 + 0) + 0.9 * (-1 + V_{\pi_b}(x)) = -1.9$$

$$V_{\pi_b}(y) = 0.2 * (-2 + 0) + 0.8 * (-2 + V_{\pi_b}(y)) = -1.8$$

Policy improvement

$$V_{\pi_b}(x) = 0.2 * (-1 + 0) + 0.8 * (-1 + V_{\pi_b}(x))$$

$$V_{\pi_b}(y) = 0.2 * (-2 + 0) + 0.8 * (-2 + V_{\pi_b}(y))$$

$$V_{\pi_b}(z) = -0.1 + (-1.8 + 0.9 V_{\pi_b}(z)) = -1.8$$

$$V_{\pi_b}(z) = -1.8 + 0.9 V_{\pi_b}(z) = -1.8$$

$$\therefore V_{\pi_b}(x) = -1.8 < V_{\pi_b}(y) = -1.8$$

$$V_{\pi_b}(y) = -1.8 > V_{\pi_b}(x) = -1.8$$

$$\pi' \neq \pi \rightarrow$$
 again Policy Evaluation

Policy Evaluation

$$V_{\pi_b}(x) = 0.2 * (-1 + 0) + 0.8 * (-1 + V_{\pi_b}(x))$$

$$V_{\pi_b}(y) = 0.2 * (-2 + 0) + 0.8 * (-2 + V_{\pi_b}(y))$$

$$V_{\pi_b}(z) = -0.1 + (-1.8 + 0.9 V_{\pi_b}(z)) = -1.8$$

$$V_{\pi_b}(z) = -1.8 + 0.9 V_{\pi_b}(z) = -1.8$$

$$\therefore V_{\pi_b}(x) = -1.8 < V_{\pi_b}(y) = -1.8$$

$$V_{\pi_b}(y) = -1.8 > V_{\pi_b}(x) = -1.8$$

$$\pi' \neq \pi \rightarrow$$
 again Policy Evaluation

Policy Evaluation

$$V_{\pi_b}(x) = 0.2 * (-1 + 0) + 0.8 * (-1 + V_{\pi_b}(x))$$

$$V_{\pi_b}(y) = 0.2 * (-2 + 0) + 0.8 * (-2 + V_{\pi_b}(y))$$

$$V_{\pi_b}(z) = -0.1 + (-1.8 + 0.9 V_{\pi_b}(z)) = -1.8$$

$$V_{\pi_b}(z) = -1.8 + 0.9 V_{\pi_b}(z) = -1.8$$

$$\therefore V_{\pi_b}(x) = -1.8 < V_{\pi_b}(y) = -1.8$$

$$V_{\pi_b}(y) = -1.8 > V_{\pi_b}(x) = -1.8$$

$$\pi' \neq \pi \rightarrow$$
 again Policy Evaluation

Policy Evaluation

$$V_{\pi_b}(x) = 0.2 * (-1 + 0) + 0.8 * (-1 + V_{\pi_b}(x))$$

$$V_{\pi_b}(y) = 0.2 * (-2 + 0) + 0.8 * (-2 + V_{\pi_b}(y))$$

$$V_{\pi_b}(z) = -0.1 + (-1.8 + 0.9 V_{\pi_b}(z)) = -1.8$$

$$V_{\pi_b}(z) = -1.8 + 0.9 V_{\pi_b}(z) = -1.8$$

$$\therefore V_{\pi_b}(x) = -1.8 < V_{\pi_b}(y) = -1.8$$

$$V_{\pi_b}(y) = -1.8 > V_{\pi_b}(x) = -1.8$$

$$\pi' \$$

Q. 4

```
=====
== Optimal State Value ==
=====

[[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]]
=====

=====
=====

== Optimal Policy ==
=====

[0, 0] = ['east']
[0, 1] = ['north', 'south', 'west', 'east']
[0, 2] = ['west']
[0, 3] = ['north', 'south', 'west', 'east']
[0, 4] = ['west']

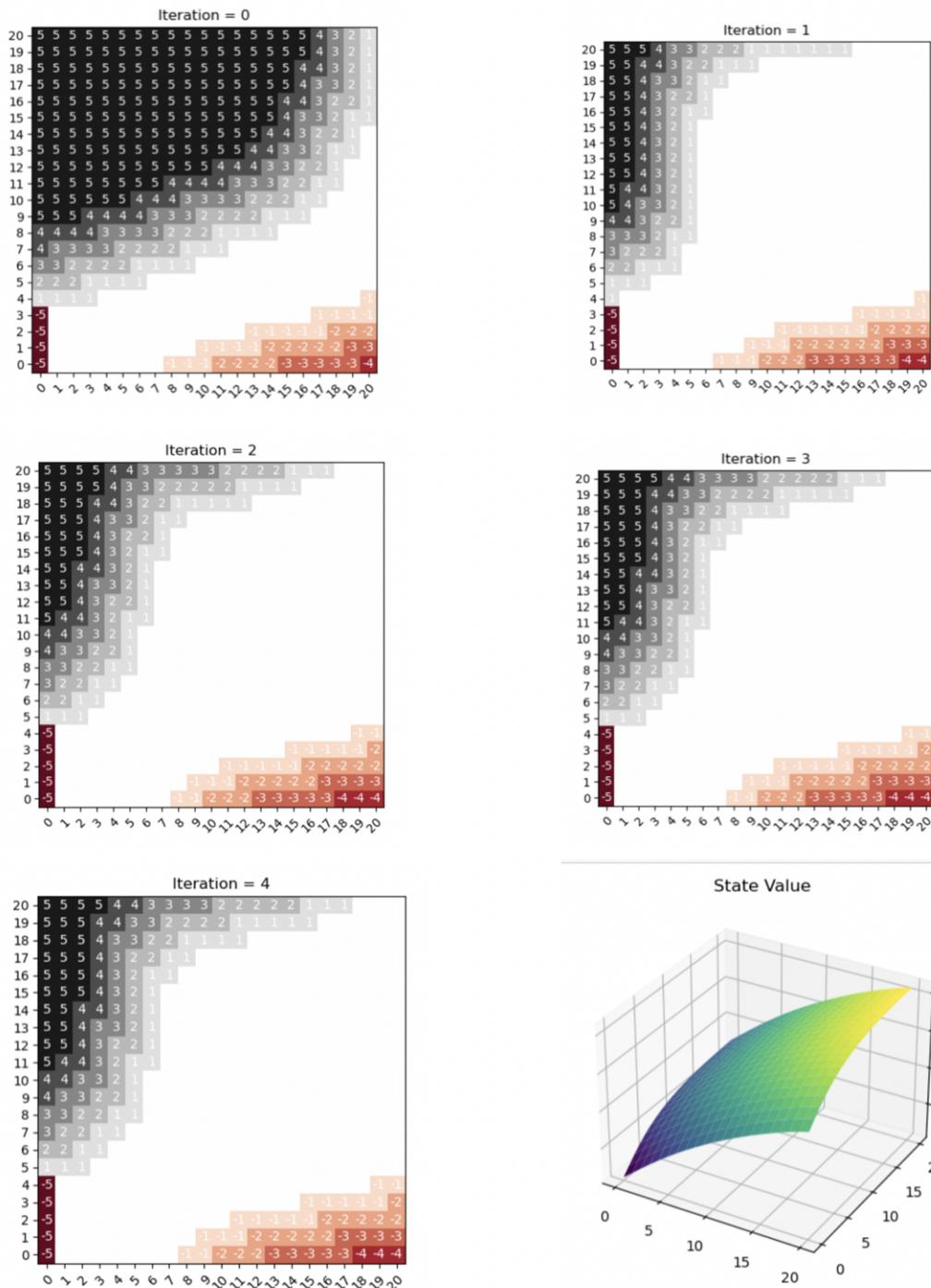
-----
[1, 0] = ['north', 'east']
[1, 1] = ['north']
[1, 2] = ['north', 'west']
[1, 3] = ['west']
[1, 4] = ['west']

-----
[2, 0] = ['north', 'east']
[2, 1] = ['north']
[2, 2] = ['north', 'west']
[2, 3] = ['north', 'west']
[2, 4] = ['north', 'west']

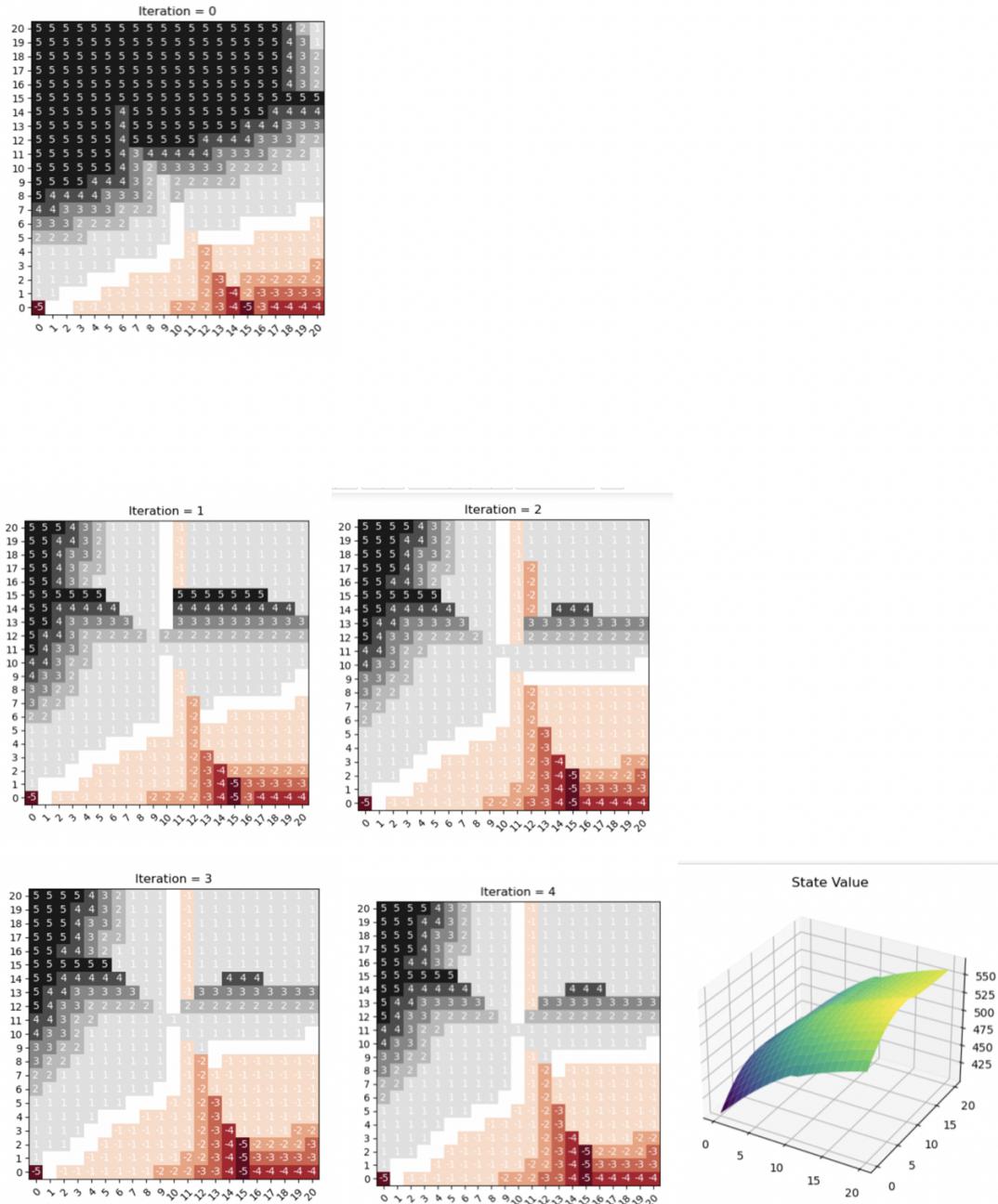
-----
[3, 0] = ['north', 'east']
[3, 1] = ['north']
[3, 2] = ['north', 'west']
[3, 3] = ['north', 'west']
[3, 4] = ['north', 'west']

-----
[4, 0] = ['north', 'east']
[4, 1] = ['north']
[4, 2] = ['north', 'west']
[4, 3] = ['north', 'west']
[4, 4] = ['north', 'west']
```

Q.5 a



Q 5b



Reward changes with modifications

- One of Jack's employees at the first location rides a bus home each night and lives near the second location. She is happy to shuttle one car to the second location for free. Each additional car still costs 2, as do all cars moved in the other direction.
- In addition, Jack has limited parking space at each location. If more than 10 cars are kept overnight at a location (after any moving of cars), then a total additional cost of 4 must be incurred to use a second parking lot (independent of how many cars are kept there). (Each location has a separate overflow lot, so if both locations have > 10 cars, the total additional cost is 8.)

} As these
are the
conditions
for the
reward
changes

$$\text{Reward} \leftarrow \begin{aligned} & \text{no_of_cars_rented} * 10 \\ & - 2 * (\text{no_of_Cars_moved} - 1, 0) \\ & - \text{Parking_charges} \end{aligned}$$

Resulting \leftarrow the location where employees ride the shuttle
the cars moves more in one direction
since cost of moving cars decreased
but because of the parking fees moving
cars becomes costly after certain limit