

Please refer to the "read_me.txt" for running the code.

Answer for q0:

(a)(b) I agree to share code. (I would argue that the shared code should be given bonus)

Answer for q1:

$$1. (a) \quad V_*(s) = \max_a Q_*(s, a)$$

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

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

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

$$(e) \quad V_*(s) = \sum_a \pi(a|s) \left(r(s, a) + \gamma \sum_{s'} p(s' | s, a) \cdot V_*(s') \right)$$

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

$$V_*(s) = \sum_a \pi(a|s) \left(r(s, a) + \sum_{s'} p(s' | s, a) \cdot \gamma \cdot V_*(s') \right)$$

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

$V_*(s)$ and $Q_*(s, a)$ could be also written as follows.

$$V_*(s) = \max_a \left(r(s, a) + \sum_{s'} p(s' | s, a) \cdot \gamma \cdot V_*(s') \right)$$

$$Q_*(s, a) = r(s, a) + \sum_{s'} p(s' | s, a) \cdot \gamma \cdot \max_{a'} Q_*(s', a')$$

Answer for q2:

I prefer the method1 in 2(a) due to that `np.argmax()` would return the first best value's index and I use this method in q5(c) to output the optimal policy.

2. (a) Method 1: give all the ~~distribution~~ probability to the first best action, in this case, the agent would always select the first best action in each state.

Method 2: Due to that the best actions have ^{the} same values, we could use best action's values to check ~~whether~~ ^{whether} the two policy are the same.

Pseudocode for Method 1:

Keep actions in $A(s)$ having fixed sequence.

$$\pi(s) = \text{first-argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

Pseudocode for Method 2:

For each $s \in S$:

$$\text{old-action} \leftarrow \pi(s)$$

$$V_1 = \sum p(s',r|s, \text{old-action}) [r + \gamma V(s')]$$

$$\pi(s) \leftarrow \text{argmax}_a \sum p(s',r|s,a) [r + \gamma V(s')]$$

$$V_2 = \sum p(s',r|s, \pi(s)) [r + \gamma V(s')]$$

if $V_1 \neq V_2$, policy-stable \leftarrow false

(b) Value functions doesn't suffer from the issue:

The best actions have the same values

Answer for q3:

3.a

1. Initialization

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

2. Policy Evaluation

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in S$:

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

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

$$Q(s,a) \leftarrow \sum p(s',r|s,a) [r + \gamma \sum_a \pi(a|s') \cdot Q(s',a)]$$

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

until $\Delta < \theta$

3. Policy Improvement

policy-stable \leftarrow true

For each $s \in S$:

$$\text{old_action} = \pi(s)$$

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

If $\text{old_action} \neq \pi(s)$, then policy-stable \leftarrow false

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

b.

$$q_{\pi}(s,a) = \sum_{s',r} p(s',r|s,a) [r + \gamma \cdot \max_{a'} q(s',a)]$$

Answer for q4:

4. (a) state x prefer^{to} take action c

state y may prefer to take action b

(b) Initial value $V(x)=0$, $V(y)=0$, $V(z)=0$
Initial policy $\pi(x)=c$, $\pi(y)=c$

policy evaluation: $\pi(x)=c$, $\pi(y)=c$

$$\begin{cases} V(x) = 0.1(-1 + V(z)) + 0.9(-1 + V(x)) \\ V(y) = 0.1(-2 + V(z)) + 0.9(-2 + V(y)) \end{cases}$$

$$\begin{cases} V(x) = -10 \\ V(y) = -20 \end{cases}$$

policy improvement:

$$Q(x, c) = -10$$

$$Q(x, b) = 0.8(-1 + V(y)) + 0.2(-1 + V(x)) = -19$$

$$\pi(x) = \operatorname{argmax}_a Q(x, a) = c$$

$$Q(y, c) = -20$$

$$Q(y, b) = 0.8(-2 + V(x)) + 0.2(-2 + V(y)) = -14$$

$$\pi(y) = \operatorname{argmax}_a Q(y, a) = b$$

Policy evaluation: $\pi(x)=c$, $\pi(y)=b$

$$\begin{cases} V(x) = 0.1(-1 + V(z)) + 0.9(-1 + V(x)) \\ V(y) = 0.2(-2 + V(y)) + 0.8(-2 + V(x)) \end{cases}$$

$$\begin{cases} V(x) = -10 \\ V(y) = -12.5 \end{cases}$$

Policy improvement:

$$\begin{aligned} Q(x, b) &= 0.2(-1 + V(x)) + 0.8(-1 + V(y)) \\ &= -13 \end{aligned}$$

$$Q(x, c) = -10$$

$$\pi(x) = \arg\max_a Q(x, a) = c$$

$$Q(y, b) = -12.5$$

$$\begin{aligned} Q(y, c) &= 0.1(-2 + V(z)) + 0.9(-2 + V(y)) \\ &= -13.25 \end{aligned}$$

$$\pi(y) = \arg\max_a Q(y, a) = b$$

Policy optimal!

The optimal action: $\pi(x)=c$; $\pi(y)=b$

The optimal value: $V(x)=-10$; $V(y)=-12.5$

4(c). Initial value $V(x)=0$, $V(y)=0$, $V(z)=0$

Initial policy $\pi(x)=b$, $\pi(y)=b$.

Policy evaluation: (Initial policy)

$$\begin{cases} V(x) = 0.2(-1 + V(x)) + 0.8(-1 + V(y)) & \textcircled{1} \\ V(y) = 0.2(-2 + V(y)) + 0.8(-2 + V(x)) & \textcircled{2} \end{cases}$$

$$\begin{cases} 0.8V(x) = -1 + 0.8V(y) \\ 0.8V(y) = -2 + 0.8V(x) \end{cases} \rightarrow \underline{0.8V(x) = -1 - 2 + 0.8V(x)} \quad \textcircled{3}$$

there is no solution for $\textcircled{3}$, and thus it couldn't converge when doing evaluation step by step.

Using discount: $0 < \gamma < 1$

policy evaluation: (Initial policy)

$$\begin{cases} V(x) = 0.2(-1 + \gamma \cdot V(x)) + 0.8(-1 + \gamma \cdot V(y)) \\ V(y) = 0.2(-2 + \gamma \cdot V(y)) + 0.8(-1 + \gamma \cdot V(x)) \end{cases}$$
$$\begin{cases} V(x) = \frac{-1 - 1.4\gamma}{(1-0.2\gamma)^2 - (0.8\gamma)^2} \\ V(y) = \frac{-2 - 0.4\gamma}{(1-0.2\gamma)^2 - (0.8\gamma)^2} \end{cases}$$

$$\therefore 0 < \gamma < 1 \rightarrow \begin{cases} (1-0.2\gamma)^2 - (0.8\gamma)^2 > 0 \\ -1 - 1.4\gamma > -2 - 0.4\gamma \end{cases} \rightarrow V(x) > V(y)$$

\therefore discount helps!

policy improvement:

$$Q(x, b) = \frac{-1 - 1.4\gamma}{(1-0.2\gamma)^2 - (0.8\gamma)^2}$$

$$Q(x, c) = 0.9(-1 + r \cdot V(x)) + 0.1(-1 + r \cdot V(z))$$

$$= -1 + 0.9r \cdot \frac{-1 - 1.4r}{(1-0.2r)^2 - (0.8r)^2}$$

$$Q(y, b) = \frac{-2 - 0.4r}{(1-0.2r)^2 - (0.8r)^2}$$

$$Q(y, c) = 0.9(-2 + r \cdot V(y)) + 0.1(-2 + r \cdot V(z))$$

$$= -2 + 0.9r \cdot \frac{-2 - 0.4r}{(1-0.2r)^2 - (0.8r)^2}$$

$$\text{For } \pi(x) = \arg \max_a Q(x, a)$$

$$\pi(y) = \arg \max_a Q(y, a)$$

and each Q value involve r . So the optimal value function depends on discount factor.

For example: If $r=0$, $Q_*(y, b) = Q_*(y, c) = -2 \rightarrow \pi_*(y) = b \text{ or } c$

If $r=1$ Without the action b in both state, $\pi_*(y) = b$.
initial policy for

due to the fact the optimal policy doesn't depend on initial policy.

We could conclude that optimal policy depend on r .

Answer for q5:

a).

run python value_estimation.py to get the value function for the equiprobable random policy in the figures below

3.314	8.793	4.431	5.326	1.496
1.526	2.996	2.253	1.911	0.55
0.055	0.742	0.676	0.361	-0.4
-0.97	-0.432	-0.352	-0.583	-1.18
-1.854	-1.342	-1.226	-1.42	-1.972

b) run value_iteration.py to get the optimal value function and the optimal policy in the figures below

21.976	24.419	21.977	19.419	17.477
19.779	21.977	19.779	17.801	16.021
17.801	19.779	17.801	16.021	14.419
16.021	17.801	16.021	14.419	12.977
14.419	16.021	14.419	12.977	11.679

right,	left,right,up,down,	left,	left,right,up,down,	left,
right,up,	up,	left,up,	left,	left,
right,up,	up,	left,up,	left,up,	left,up,
right,up,	up,	left,up,	left,up,	left,up,
right,up,	up,	left,up,	left,up,	left,up,

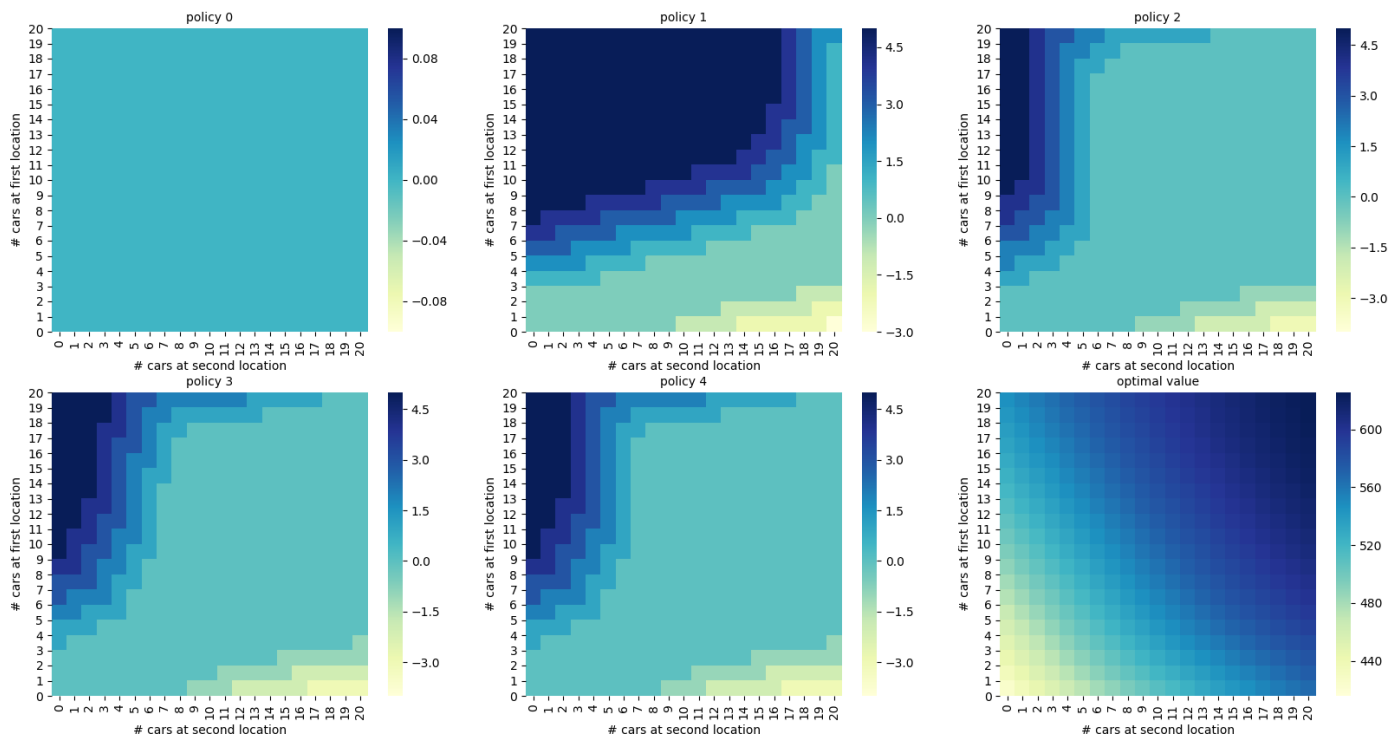
c) run policy iteration.py to get the optimal value function and the optimal policy in the figures below

21.977	24.419	21.977	19.419	17.477
19.78	21.977	19.78	17.802	16.022
17.802	19.78	17.802	16.022	14.419
16.022	17.802	16.022	14.419	12.977
14.419	16.022	14.419	12.977	11.68

right	right	left	right	left
right	up	left	left	left
right	up	left	left	left
right	up	left	left	left
right	up	left	left	left

Answer for q6:

a) Run python jack_car_a.py --constant_return to get the optimal value function and the optimal policy in the figures below



Some details:

Constant return is used to simplify the model: (You could also run `python jack_car_a.py` to get the accurate model)

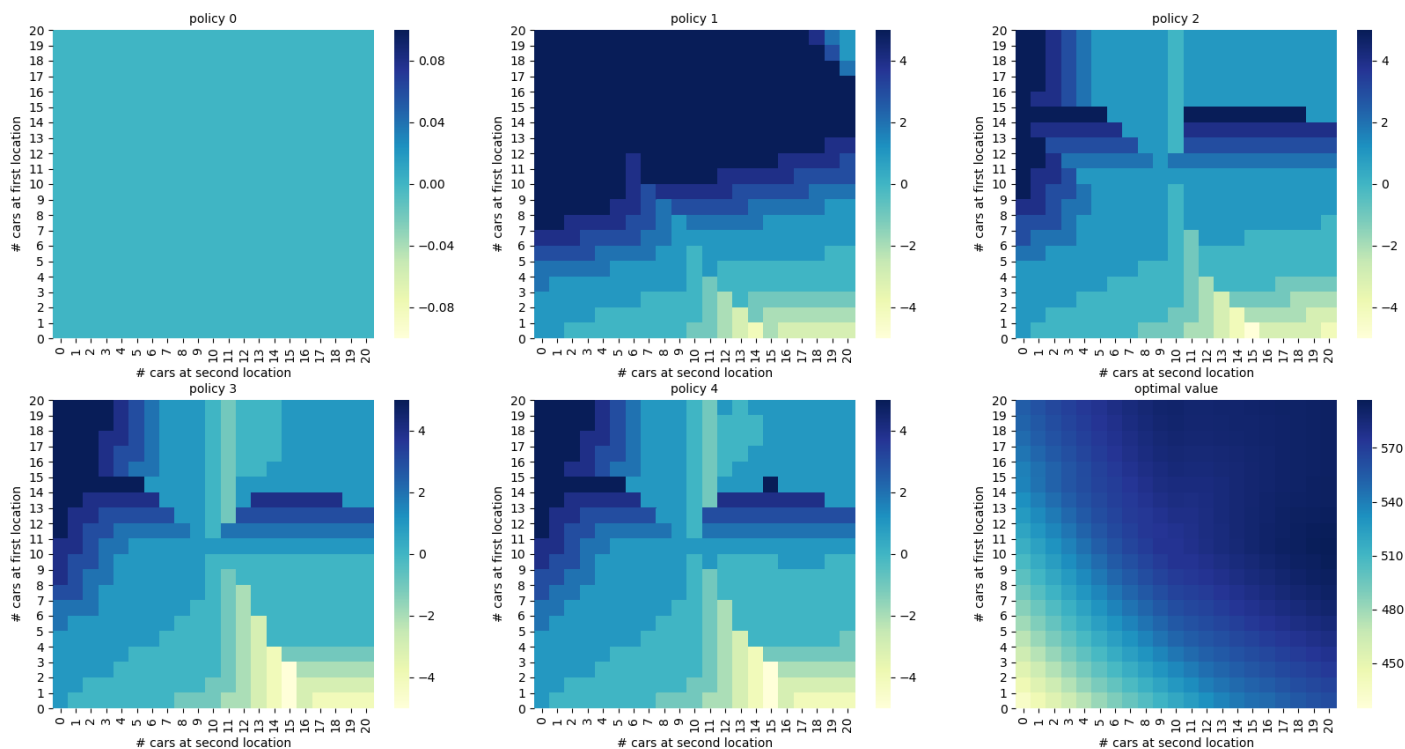
In order to check the output quickly, constant returns instead of the random returns from poisson distribution would be more efficient. Specifically, the expectation of poisson distribution is λ . In this case, the expectation of returns in the first location is 3 and the expectation returns in the second location is 2. It would leave the optimal policy and value almost the same.

b):

dynamic function changes:

- the employee's help could save 2 cost if the action is moving one or more cars from the first location to the second location and thus increase reward compared with the previous dynamic function
- the parking would cost 4 for each location if the parking number bigger than 10 and thus would decrease the reward compared with the previous dynamic function

run `jack_rental_car_b.py --constant_return` to get the optimal value function and the optimal policy in the figures



Details:

You could also run `python jack_car_b.py` to use the accurate the model:

In order to check the output quickly, constant returns instead of the random returns from poisson distribution would be more efficient. Specifically, the expectation of poisson distribution is λ . In this case, the expectation of returns in the first location is 3 and the expectation returns in the second location is 2. It would leave the optimal policy and value almost the same.

Policy difference:

- States that have 10 more cars in location1 would like to move more cars to location 2 to avoid parking fee and use the employee's help
- States that have 10 more cars in one location and less than 10 cars in the another location would like to move cars (won't excess 2) to the location that has less cars to avoid parking fee.

Answer for q7:

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

given any function f, g :

$$\begin{aligned} |f(a) - g(a)| &\geq f(a) - g(a) \\ |f(a) - g(a)| + g(a) &\geq f(a) \quad (1) \end{aligned}$$

$$\max (|f(a) - g(a)| + g(a)) \geq \max f(a) \quad (2)$$

$$\therefore \max |f(a) - g(a)| + \max g(a) \geq \max (|f(a) - g(a)| + g(a))$$

$$\therefore \max |f(a) - g(a)| + \max g(a) \geq \max f(a) \quad (3)$$

$$\max |f(a) - g(a)| \geq \max f(a) - \max g(a) \quad (4)$$

It's easy to use the symmetry trick to replace $f(a), g(a)$:

$$\max |f(a) - g(a)| \geq \max g(a) - \max f(a) \quad (5)$$

From (4), (5),

$$\max |f(a) - g(a)| \geq |\max f(a) - \max g(a)|$$

$$7. (b) V_i = \begin{bmatrix} V_i(s_1) \\ V_i(s_2) \\ \vdots \\ V_i(s_n) \end{bmatrix} \quad V_i' = \begin{bmatrix} V_i'(s_1) \\ V_i'(s_2) \\ \vdots \\ V_i'(s_n) \end{bmatrix}$$

$$\|BV_i - BV_i'\| = \begin{bmatrix} \max_a \sum_{s_1, r} p(s_1', r | s_1, a) (r + V_i(s_1)) - \max_a \sum_{s_1, r} p(1) (r + V_i'(s_1)) \\ \max_a \sum_{s_2, r} p(s_2', r | s_2, a) (r + V_i(s_2)) - \max_a \sum_{s_2, r} p(1) (r + V_i'(s_2)) \\ \vdots \\ \max_a \sum_{s_n, r} p(s_n', r | s_n, a) (r + V_i(s_n)) - \max_a \sum_{s_n, r} p(1) (r + V_i'(s_n)) \end{bmatrix}$$

$$\leq \begin{bmatrix} \max_a \left| \sum_{s_1, r} p(s_1', r | s_1, a) \gamma (V_i(s_1) - V_i'(s_1)) \right| \\ \max_a \left| \sum_{s_2, r} p(s_2', r | s_2, a) \cdot \gamma (V_i(s_2) - V_i'(s_2)) \right| \\ \vdots \\ \max_a \left| \sum_{s_n, r} p(s_n', r | s_n, a) \cdot \gamma (V_i(s_n) - V_i'(s_n)) \right| \end{bmatrix} \quad (2)$$

From (2), we could get

$$\|BV_i - BV_i'\|_\infty \leq \left\| \begin{bmatrix} \max_a \left| \sum_{s_1, r} p(s_1', r | s_1, a) \gamma (V_i(s_1) - V_i'(s_1)) \right| \\ \max_a \left| \sum_{s_2, r} p(s_2', r | s_2, a) \cdot \gamma (V_i(s_2) - V_i'(s_2)) \right| \\ \vdots \\ \max_a \left| \sum_{s_n, r} p(s_n', r | s_n, a) \cdot \gamma \cdot (V_i(s_n) - V_i'(s_n)) \right| \end{bmatrix} \right\|_\infty$$

(3)

$$\therefore \sum p(s_i', r | s_i, a) \leq 1$$

$$\therefore (3) \leq \gamma \|V_i - V_i'\|_\infty$$

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

7 (b) Method 1:

Assume that $\{V_0, V_1, V_2, V_3 \dots V_n\}$ is sequence of values.

$$\underline{V_n = BV_{n-1}}, \text{ and } \Delta = \|V_0 - V_1\|_\infty :$$

from 7(b), Bellman backup operator is a ~~backup~~ contraction mapping.

$$\|BV_0 - BV_1\|_\infty \leq \gamma \|V_0 - V_1\|_\infty = \gamma \cdot \Delta$$

$$\therefore \|BV_n - BV_{n+1}\|_\infty \leq \gamma^{n+1} \|V_0 - V_1\|_\infty = \gamma^{n+1} \cdot \Delta$$

$\lim_{n \rightarrow \infty} \|BV_n - BV_{n+1}\|_\infty \leq \lim_{n \rightarrow \infty} \gamma^{n+1} \cdot \Delta$, where $0 < \gamma < 1$, Δ is a constant.

$$\therefore \left. \begin{array}{l} \lim_{n \rightarrow \infty} \|BV_n - BV_{n+1}\|_\infty = 0 \quad ① \\ V_n = BV_{n-1} \rightarrow BV_n = V_{n+1} \quad ② \end{array} \right\} \rightarrow \lim_{n \rightarrow \infty} BV_n = V_n$$

Method 2: use Cauchy sequence, metric space to prove the fixed points

(c) A is the sequence contains x_0, x_1, \dots, x_n , where $x_n = \beta(x_{n-1})$

d is a metric function: $d(x_n, x_m) = \|x_n - x_m\|_\infty$

$d(x, y) \leq d(x, z) + d(y, z)$ triangle inequality

In this case (A, d) is a metric space ①

From b, we could get

$$\|Bx_1 - Bx_0\|_\infty \leq r \|x_1 - x_0\|_\infty \rightarrow d(x_2, x_1) \leq r d(x_1, x_0) \quad ②$$

It's easy to find $d(x_{n+1}, x_n) \leq q^n d(x_1, x_0)$, where $q = r$

For Bellman operator is a contraction mapping, we can show $\{x_n\}$ is a Cauchy sequence.

let $m, n \in \mathbb{N}$ and $m > n$,

$$\begin{aligned} d(x_m, x_n) &\leq d(x_m, x_{m-1}) + d(x_{m-1}, x_{m-2}) + \dots + d(x_{n+1}, x_n) \quad (\text{triangle inequality}) \\ &\leq q^{m-1} d(x_1, x_0) + q^{m-2} d(x_1, x_0) + \dots + q^n d(x_1, x_0) \\ &= q^n d(x_1, x_0) \sum_{k=0}^{m-n-1} q^k \\ &\leq q^n d(x_1, x_0) \sum_{k=0}^{\infty} q^k = \underbrace{q^n d(x_1, x_0) \left(\frac{1}{1-q} \right)}_{\Delta} \quad ③ \end{aligned}$$

let $\varepsilon > 0$ be arbitrary, since $q = r \in [0, 1)$, there would be a large $C \in \mathbb{N}$,

$$q^C < \frac{\varepsilon(1-q)}{d(x_1, x_0)} \quad ④$$

If m, n larger than C

$$\begin{aligned} \therefore d(x_m, x_n) &\leq q^n d(x_1, x_0) \left(\frac{1}{1-q} \right) \xrightarrow{q^n < q^C} q^C < \frac{\varepsilon(1-q)}{d(x_1, x_0)} \quad ⑤ \\ &< \frac{\varepsilon(1-q)}{d(x_1, x_0)} \cdot d(x_1, x_0) \left(\frac{1}{1-q} \right) = \varepsilon \end{aligned}$$

So for every positive real number ϵ , there is a positive integer N

Such that for all natural numbers $m, n > N$

$d(X_m - X_n) < \epsilon$, the sequence is Cauchy sequence.

In this case, the sequence $\{X_n\}$ has a limit $X^* \in X$

$$X^* = \lim_{n \rightarrow \infty} X_n = \lim_{n \rightarrow \infty} B(X_{n-1}) = B(\lim_{n \rightarrow \infty} X_{n-1}) = B(X^*)$$

, which shows value iteration ~~converges~~ converges to a fix point. (converge)

If there are two fixed points P_1 and P_2 :

$$\|B(P_1) - B(P_2)\|_{L-\infty} \leq \gamma \cdot \|P_1 - P_2\|_{L-\infty}$$

$$\because B(P_1) = P_1, \quad B(P_2) = P_2 \quad (\text{fixed points' property})$$

$$\therefore \|P_1 - P_2\|_{L-\infty} \leq \gamma \cdot \|P_1 - P_2\|_{L-\infty} \quad \textcircled{1}$$

$$\therefore 0 < \gamma < 1$$

\therefore equation ① would be correct if only $P_1 = P_2$

, which shows the fixed point is unique. (unique)

$$\text{When } B(X^*) = X^*, \quad X^*(s) = \max_{\alpha} \sum_{s', r} p(s', r | s, \alpha) [r + \gamma X^*(s')]$$

$$= \max_{\alpha} E_{x_t} [R_{t+1} + \gamma X^*(S_{t+1}) | S_t = s, A_t = \alpha]$$

$$= \max_{\alpha} E_{x_t} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = \alpha]$$

$$= \max_{\alpha} E_{x_t} [G_t | S_t = s, A_t = \alpha]$$

$$X^*(s) = \max_a E_{\pi^*} [G_t \mid s_t=s, A_t=a]$$

$$= \max_a \underline{q_{\pi^*}(s,a)}$$

which shows $X^*(s)$ is the max q_{π^*} value and thus ~~was~~ is optimal value. (optimal value function)

Reference:

<https://github.com/ShangtongZhang/reinforcement-learning-an-introduction>

https://en.wikipedia.org/wiki/Banach_fixed-point_theorem