

Q0: (a) Yes. Either anonymized or not is okay

(b) Same for EX0 and EX1

Q1:

(a)

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

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

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

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

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

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

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

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

Q2:

(a). The policy keeps switching between two optimal ones because the original pseudo code use actions as criteria to judge stability of the policy, while actually we should stop the iteration when we know we already have an optimal one

fix: change the second line from the last to:

if  $\sum_{s', r} p(s', r | s, \text{old\_action}) (r + \gamma V(s')) \neq \sum_{s', r} p(s', r | s, \pi(s)) (r + \gamma V(s'))$ , policy-stable  $\leftarrow$  false

(b). There isn't, because value iteration derives the optimal  $V^*$  first, and then calculate the policy, so there isn't switching action.

Q3:

(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  $s \in S$

$a \leftarrow \pi(s)$

$q \leftarrow Q(s, a)$

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

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

until  $\Delta < \theta$  (a small value for accuracy)

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  $\sum_{s', r} P(s', r | s, \pi(s)) [r + \gamma \sum_{a'} \pi(a' | s') q(a', s')] \neq \sum_{s', r} P(s', r | s, \text{old-action}) [r + \gamma \sum_{a'} \pi(a' | s') q(a', s')]$ ,  
then policy-stable  $\leftarrow$  false

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

(b) init: threshold  $\theta > 0$ ;  $Q(s, a)$ , for all  $s \in S^+$ , arbitrarily  $Q(s, a) = 0$  when  $s$  is terminal.

Loop:

$\Delta \leftarrow 0$

Loop for each  $s \in S$ :

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

$q \leftarrow Q(s, a)$

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

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

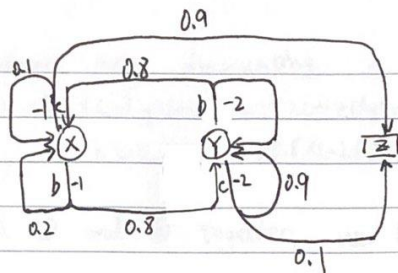
until  $\Delta < \theta$ .

Output deterministic policy,  $\pi \approx \pi^*$ , such that

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

Q4:

(a)



Since no matter what action we choose at  $y$ , it will always be worse than in  $x$ , so qualitatively optimal policy for  $x$  should be  $C$  because by doing this it can avoid entering  $y$ , which is worse. For  $y$ , it's not very clear to me, since whether entering  $x$  or  $c$  it's always a improvement, and I think whether  $b$  or  $c$  is optimal for  $y$  is depended on the discounting factor.

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

$$\begin{cases} V_{\pi}(x) = 0.1 \times (-1) + 0.9 \times (-1 + V_{\pi}(x)) \\ V_{\pi}(y) = 0.1 \times (-2) + 0.9 \times (-2 + V_{\pi}(y)) \end{cases} \Rightarrow \begin{cases} V_{\pi}(x) = -10 \\ V_{\pi}(y) = -20 \end{cases} \quad \text{evaluation 1}$$

$$\pi(s) = \arg \max_a \sum_{s',r} p(s',r|s,a)(r + V_{\pi}(s'))$$

$$V_{\pi b}(x) = 0.2 \times (-1 + V_{\pi}(x)) + 0.8 \times (-1 + V_{\pi}(y)) = -1 - 2 - 16 = -19$$

$$V_{\pi b}(y) = 0.2 \times (-2 + V_{\pi}(y)) + 0.8 \times (-2 + V_{\pi}(y)) = -2 - 4 - 8 = -14$$

$$\therefore \begin{cases} V_{\pi c}(x) = -10 \\ V_{\pi c}(y) = -20 \end{cases}$$

$$\therefore \pi'(x) = c, \pi'(y) = b.$$

improvement 1

$$\begin{cases} V_{\pi}(x) = 0.1 \times (-1) + 0.9 \times (-1 + V_{\pi}(x)) \\ V_{\pi}(y) = 0.2 \times (-2 + V_{\pi}(y)) + 0.8 \times (-2 + V_{\pi}(x)) \end{cases} \Rightarrow \begin{cases} V_{\pi}(x) = -10 \\ V_{\pi}(y) = -12.5 \end{cases} \quad \text{evaluation 2.}$$

$$V_{\pi b}(x) = 0.2 \times (-1 + V_{\pi}(x)) + 0.8 \times (-1 + V_{\pi}(y)) = -13$$

$$V_{\pi c}(y) = 0.1 \times (-2) + 0.9 \times (-2 + V_{\pi}(y)) = -13.25$$

$$\therefore \begin{cases} V_{\pi c}(x) = -10 \\ V_{\pi b}(y) = -12.5 \end{cases}$$

improvement 2.

$$\therefore \pi''(x) = c = \pi'(x), \pi''(y) = b = \pi'(y)$$

$$\therefore \text{Stabilized, } \pi_*(x) = c, \pi_*(y) = b.$$



(C) when there isn't discounting

$$\begin{cases} V_{\pi}(x) = 0.2 \times (-1 + V_{\pi}(x)) + 0.8 \times (-1 + V_{\pi}(y)) & ① \\ V_{\pi}(y) = 0.2 \times (-2 + V_{\pi}(y)) + 0.8 \times (-2 + V_{\pi}(x)) & ② \end{cases}$$

add ① and ② together, we have  $0 = -3$ , which is obviously wrong.

After adding discounting factor  $\gamma$ :

$$\begin{cases} V_{\pi}(x) = 0.2 \times (-1 + \gamma V_{\pi}(x)) + 0.8 \times (-1 + \gamma V_{\pi}(y)) \\ V_{\pi}(y) = 0.2 \times (-2 + \gamma V_{\pi}(y)) + 0.8 \times (-2 + \gamma V_{\pi}(x)) \end{cases}$$

$$\Rightarrow \begin{cases} V_{\pi}(y) = \frac{-2 - 0.4\gamma}{1 - 0.6\gamma^2 - 0.4\gamma} & = V_{\pi b}(y) \\ V_{\pi}(x) = \frac{-1 + 0.28\gamma^2 - 1.2\gamma}{(1 - 0.6\gamma^2 - 0.4\gamma)(1 + 0.2\gamma)} & = V_{\pi b}(x) \end{cases}$$

$$\begin{cases} V_{\pi c}(x) = 0.1 \times (-1) + 0.9 \times (-1 + V_{\pi}(x)) \\ \quad = -1 + 0.9 V_{\pi}(x) \end{cases} \rightarrow \text{if } V_{\pi}(x) < -10, \text{ then } \pi(x) = C$$

$$\begin{cases} V_{\pi c}(y) = -2 + 0.9 V_{\pi}(y) \end{cases} \rightarrow \text{if } V_{\pi}(y) < -20, \text{ then } \pi(y) = C.$$

By observing  $V_{\pi b}(y)$ , we can find that  $(-\infty, -2]$  is part of its value range, because when  $\gamma \rightarrow 1$ ,  $V_{\pi}(y) \rightarrow -\frac{2}{5}$ , and when  $\gamma = 0$ ,  $V_{\pi}(y) = -2$ , and it's continuous on  $\gamma \in (0, 1)$ , so we can definitely control the value of  $\gamma$  to make it higher or lower than  $-20$ , and thus control the optimal policy for  $y$ . Similarly, we can do the same thing on  $x$ .

In conclusion, optimal policy is depended on the discount factor.

Q5:

(a)

```
[[ 3.31  8.79  4.43  5.32  1.49]
 [ 1.52  2.99  2.25  1.91  0.55]
 [ 0.05  0.74  0.67  0.36 -0.4 ]
 [-0.97 -0.43 -0.35 -0.58 -1.18]
 [-1.86 -1.34 -1.23 -1.42 -1.97]]
```

Same as Figure 3.2

(b)

```
[[21.98 24.42 21.98 19.42 17.48]
 [19.78 21.98 19.78 17.8  16.02]
 [17.8  19.78 17.8  16.02 14.42]
 [16.02 17.8  16.02 14.42 12.98]
 [14.42 16.02 14.42 12.98 11.68]]
{'[0, 0]': ['r'], '[0, 1]': ['u', 'd', 'l', 'r'], '[0, 2]': ['l'],
 '[0, 3]': ['u', 'd', 'l', 'r'], '[0, 4]': ['l'], '[1, 0]': ['u',
 'r'], '[1, 1]': ['u'], '[1, 2]': ['u', 'l'], '[1, 3]': ['l'], '[1,
 4]': ['l'], '[2, 0]': ['u', 'r'], '[2, 1]': ['u'], '[2, 2]': ['u',
 'l'], '[2, 3]': ['u', 'l'], '[2, 4]': ['u', 'l'], '[3, 0]': ['u',
 'r'], '[3, 1]': ['u'], '[3, 2]': ['u', 'l'], '[3, 3]': ['u', 'l'],
 '[3, 4]': ['u', 'l'], '[4, 0]': ['u', 'r'], '[4, 1]': ['u'], '[4,
 2]': ['u', 'l'], '[4, 3]': ['u', 'l'], '[4, 4]': ['u', 'l']}
```

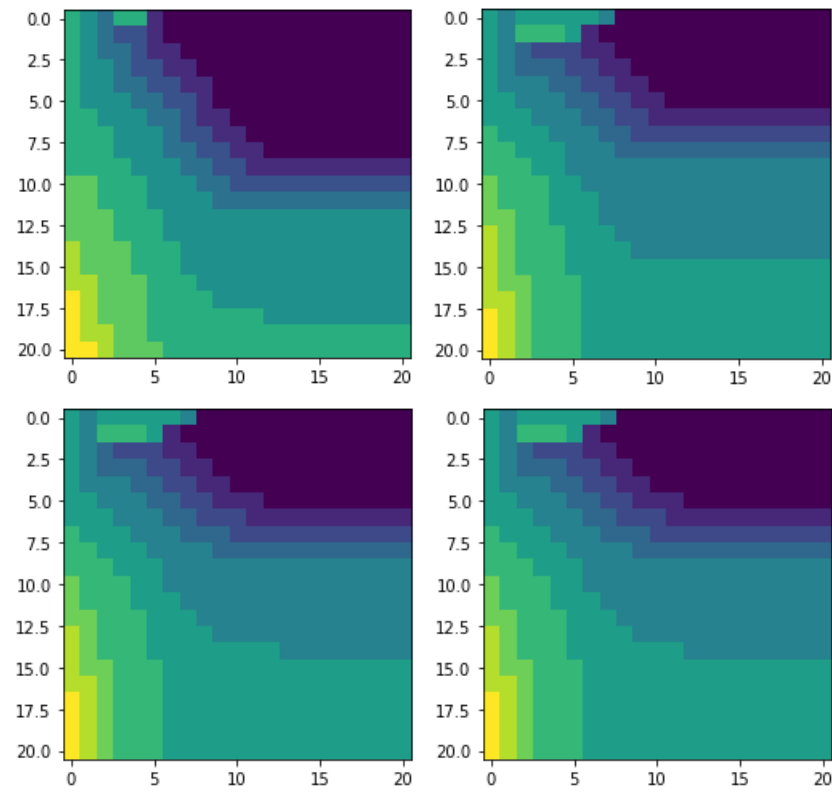
Same as Figure 3.5. The dictionary contains derived optimal policy, which key is state and value is action.

(c)

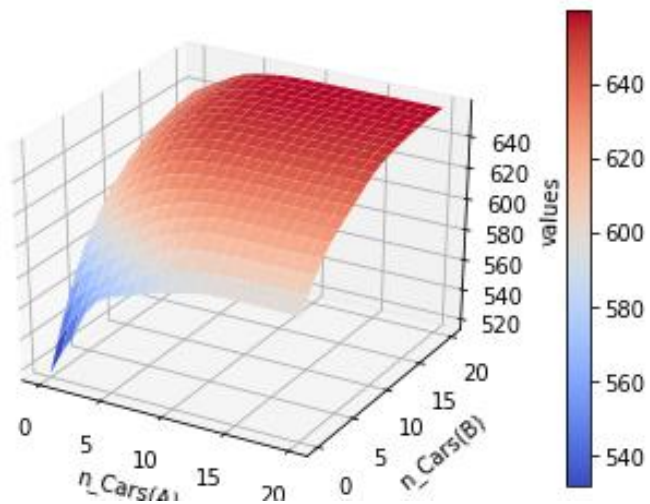
```
[[21.98 24.42 21.98 19.42 17.48]
 [19.78 21.98 19.78 17.8  16.02]
 [17.8  19.78 17.8  16.02 14.42]
 [16.02 17.8  16.02 14.42 12.98]
 [14.42 16.02 14.42 12.98 11.68]]
{'[0, 0]': ['r'], '[0, 1]': ['u', 'd', 'l', 'r'], '[0, 2]': ['l'],
 '[0, 3]': ['u', 'd', 'l', 'r'], '[0, 4]': ['l'], '[1, 0]': ['u',
 'r'], '[1, 1]': ['u'], '[1, 2]': ['u', 'l'], '[1, 3]': ['l'], '[1,
 4]': ['l'], '[2, 0]': ['u', 'r'], '[2, 1]': ['u'], '[2, 2]': ['u',
 'l'], '[2, 3]': ['u', 'l'], '[2, 4]': ['u', 'l'], '[3, 0]': ['u',
 'r'], '[3, 1]': ['u'], '[3, 2]': ['u', 'l'], '[3, 3]': ['u', 'l'],
 '[3, 4]': ['u', 'l'], '[4, 0]': ['u', 'r'], '[4, 1]': ['u'], '[4,
 2]': ['u', 'l'], '[4, 3]': ['u', 'l'], '[4, 4]': ['u', 'l']}
```

Q6:

- (a) Yellow indicates positive, blue indicates negative. The color in the middle of these four graphs indicates '0'.



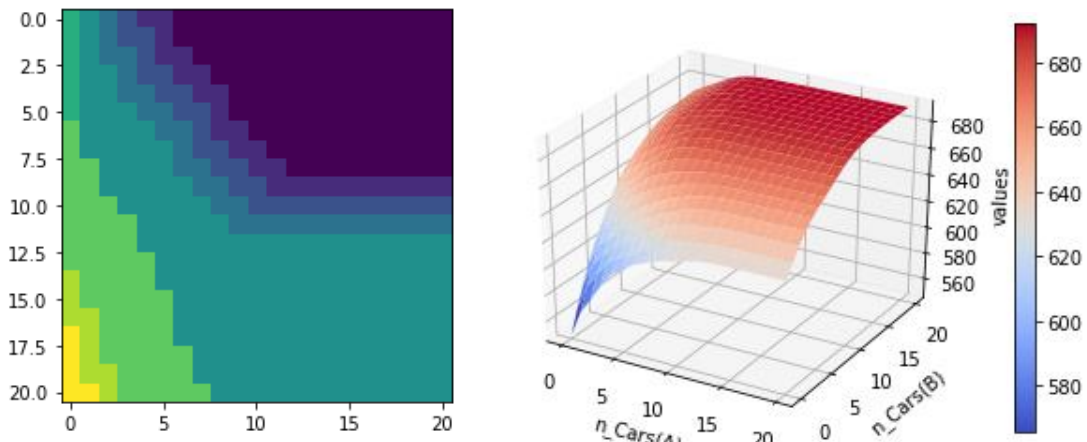
$\pi_1$  to  $\pi_4$  (top-left to bottom-right), with initial  $\pi_0$  has 0 as action for every state



3D-plot of  $v_{\pi_4}$ , maximum state value located at (20,20)

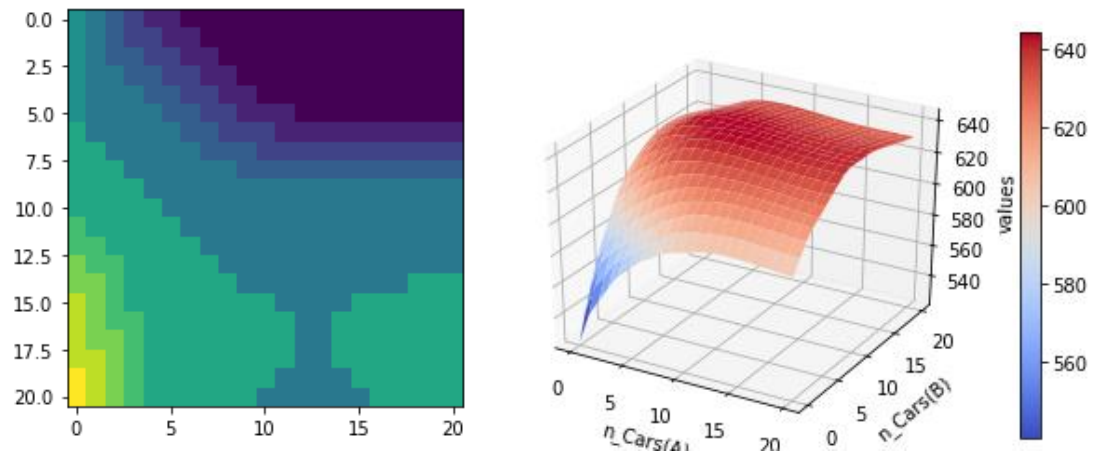
(b) The reward will change when making those modifications. When an employee is able to help switching cars the cost on switching in one of the directions will decrease by one, so I assume that employee lives at the second place and thus changed the reward when moves cars from the first place to the second from  $\text{reward} = n_{\text{rented}}*10 - 2*\text{num\_movedCars}$  to  $\text{reward} = n_{\text{rented}}*10 - 2*\max((\text{num\_movedCars}-1),0)$ , and reward function for the other direction remains unchanged. In this case the optimal policy should be tend to switch more cars in one direction.

For the second case, when the parking lot charge for extra fees there will be an extra cost which result in lower reward, so I subtracted an extra term  $\text{parking\_fee}$  if the total number of cars in that place is more than 10, as  $n_{\text{rented}}*10 - 2*\text{num\_movedCars} - \text{parking\_fee}$ . And I think the optimal policy should tend to keep less cars in both locations.



$\pi_4$  and 3D-plot of  $v_{\pi_4}$  of employee problem

As a result, the policy tends to move more cars in one direction and the other direction remains the same. The state values are also rise in general, since the cost of moving cars decreased, while the peak value still appears in (20,20).



$\pi_4$  and 3D-plot of  $v_{\pi_4}$  of parking fee problem

As a result, the policy tends to keep less cars in both locations, and sometimes it doesn't even move the cars when one of the locations has more than 15 cars because it wants to prevent being charged for extra parking fee for the other parking lot. The state values are also decreased in general, since the reward has been subtracted by extra parking fee, and it also impact the peak value which appears to be (10,10). It makes sense because keeping 10 cars in both parking lot is the best state that maximizing number of cars and avoiding extra parking fee.



Q7.

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

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

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

$\therefore \max_a f(a) - \max_a g(a) \leq \max_a f(a) - g(x)$ , for all  $x$ ,

$$\therefore \max_a f(a) - \max_a g(a) \leq \max_a f(a) - g(a) \quad (1)$$

Say,  $a_1 = \arg \max_a f(a)$ ,  $a_2 = \arg \max_a (f(a) - g(a))$

then  $f(a_1) - g(a_1) \leq f(a_2) - g(a_2)$  (2) would be easily proved.

(2)  $\Rightarrow$  (1): where  $f(a_1) - g(a_1) = \max_a f(a) - g(a)$ ,

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

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

$\therefore$  proved.

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

substitute  $f(a) \leftarrow g(a)$ , and  $g(a) \leftarrow f(a)$ , and we can prove it in the same way.

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

$$(b). |BV_i(s) - BV_i'(s)| = \left| \max_a \sum_{s',r} p(s',r|s,a) [r + V_i(s')] - \max_a \sum_{s',r} p(s',r|s,a) [r + V_i'(s')] \right| \quad (1)$$

Use equation proved from (a)

$$(1) \leq \max_a \left| \sum_{s',r} p(s',r|s,a) [r + V_i(s')] - \sum_{s',r} p(s',r|s,a) [r + V_i'(s')] \right|$$

$$= \max_a \left| \sum_{s',r} p(s',r|s,a) (r + V_i(s') - r - V_i'(s')) \right|$$

$$= \max_a \left| \sum_{s',r} p(s',r|s,a) \cdot \delta(V_i(s') - V_i'(s')) \right| \quad \text{hold for all } s \in S$$

for  $n$ -length Vector  $V_i$  and  $V_i'$ :

$$\|BV_i - BV_i'\| = \max \left\{ \begin{array}{c} |BV_i(s_1) - BV_i'(s_1)| \\ |BV_i(s_2) - BV_i'(s_2)| \\ \vdots \\ |BV_i(s_n) - BV_i'(s_n)| \end{array} \right\} \leq \max \left\{ \begin{array}{c} \max_a \left| \sum_{s',r} p(s',r|s,a) \delta(V_i(s_1) - V_i'(s_1)) \right| \\ \max_a \left| \sum_{s',r} p(s',r|s,a) \delta(V_i(s_2) - V_i'(s_2)) \right| \\ \vdots \\ \max_a \left| \sum_{s',r} p(s',r|s,a) \delta(V_i(s_n) - V_i'(s_n)) \right| \end{array} \right\} \quad (2)$$

$$② \leq \max \left\{ \begin{array}{c} \max_a \left| \sum_{s,r} p(s,r|s,a) \cdot \delta \|V_i - V_i'\| \right| \\ \vdots \\ \text{duplicate} \cdot \vdots \cdot x_n \\ \vdots \\ \max_a \left| \sum_{s,r} p(s,r|s,a) \cdot \delta \|V_i - V_i'\| \right| \end{array} \right\}$$

$$= \max \left\{ \begin{array}{c} \delta \|V_i - V_i'\| \\ \delta \|V_i - V_i'\| \\ \vdots \\ \delta \|V_i - V_i'\| \end{array} \right\} = \delta \|V_k - V_k'\|.$$

$\therefore$  proved.

(C) Banach fixed point theorem:  $d(T(x), T(y)) \leq q d(x, y)$

in our case:  $d(x, y) = \|x - y\|$ ,

$T(x) = \beta x$ ,  $\beta$  is contracting mapping

$q = \delta$

use equation from (b):  $\|\beta V_i - \beta V_i'\|_\infty \leq \delta \|V_i - V_i'\|_\infty$  (3)

proof:

first for any positive integer  $n$ :  $\|V_m - V_n\| \leq \delta^n \|V_1 - V_0\|$  (use 3).

then, for any  $m, n$ :

$$\begin{aligned} \|V_m - V_n\| &\leq \|V_m - V_{m-1}\| + \|V_{m-1} - V_{m-2}\| + \dots + \|V_{n+1} - V_n\| \\ &\leq \delta^{m-1} \|V_1 - V_0\| + \delta^{m-2} \|V_1 - V_0\| + \dots + \delta^n \|V_1 - V_0\| \\ &= \delta^n \|V_1 - V_0\| \sum_{k=0}^{m-n-1} \delta^k \\ &\leq \delta^n \|V_1 - V_0\| \sum_{k=0}^{\infty} \delta^k \\ &= \delta^n \|V_1 - V_0\| \frac{1}{1-\delta} \end{aligned}$$

We can find a large  $N$  that  $\delta^N < \frac{\epsilon(1-\delta)}{\|V_1 - V_0\|}$ ,  $\epsilon > 0$  is arbitrary value

$$\therefore \|V_m - V_n\| \leq \delta^n \|V_1 - V_0\| \frac{1}{1-\delta} < \epsilon$$

$\therefore$  the sequence  $V_i$  is Cauchy, and thus must have a fix point.

$$V^* = \lim_{n \rightarrow \infty} V_n = \lim_{n \rightarrow \infty} \beta V_{n-1} = \beta \left( \lim_{n \rightarrow \infty} V_{n-1} \right) = \beta V^*$$

$\therefore$  it converges to a fix point.

Unique: if it has two fix point,  $V_1, V_2$  then.

$$V_1 = \beta V_1, V_2 = \beta V_2.$$

$$\|\beta V_1 - \beta V_2\|_\infty = \|\beta(V_1 - V_2)\|_\infty \leq \beta \|V_1 - V_2\|_\infty$$

$$\because \beta \in (0, 1)$$

$$\therefore \|V_1 - V_2\| = 0$$

$\therefore$  unique fix point.

according to eq. 3.19 from the book,

$$\beta V_*(s) = \max_a E[R_{t+1} + \beta V_*(S_{t+1}) | S_t = s, A_t = a]$$

$$= \max_{a \in A} q_{\pi_*}(s, a)$$

$$= V_*(s)$$

so this unique fixed point satisfied Bellman optimal equation.

or equivalent to