

# Reinforcement Learning

AI 3000

Assignment No: 01

K. Surya Prakash

EE18BTEUH 11026

p1)  
@

Markov Reward Process

States:  $\langle S_0, S_1, S_2, S_3, S_4 \rangle$

$S_0$ : Start State

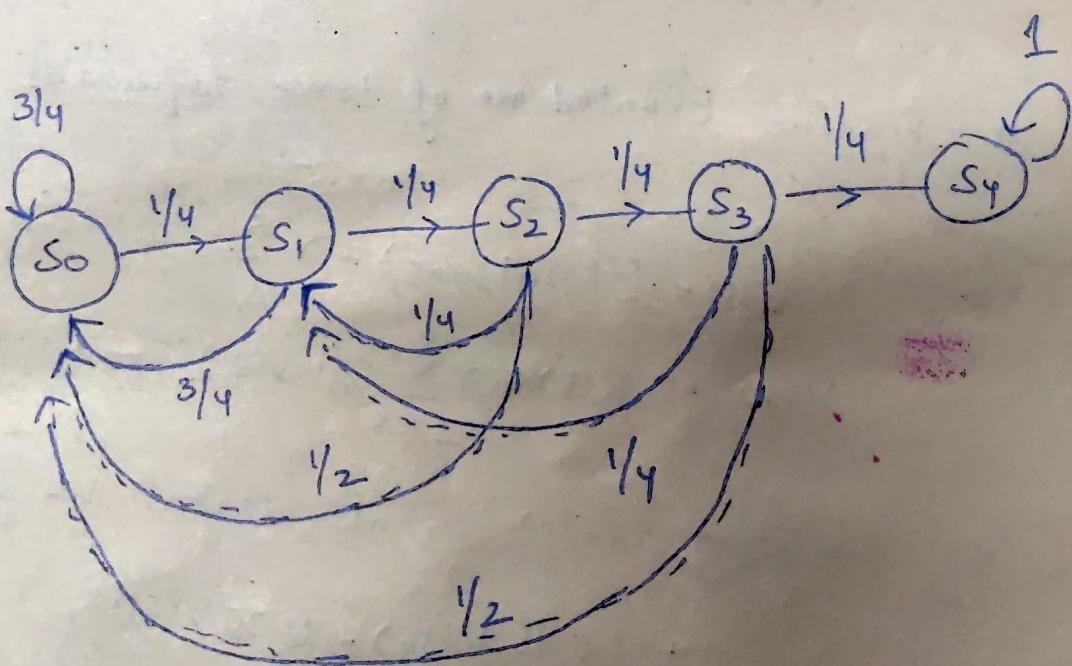
$S_1$ : When a 1 appears on the die

$S_2$ : When 2 appears after 1

$S_3$ : When 3 appears after 2

$S_4$ : When 4 appears after 3.

State transition Diagram:



	$s_0$	$s_1$	$s_2$	$s_3$	$s_4$
$s_0$	0.75	0.25	0	0	0
$s_1$	0.5	0.25	0.25	0	0
$s_2$	0.5	0.25	0	0.25	0
$s_3$	0.5	0.25	0	0	0.25
$s_4$	0	0	0	0	1

#  $s_4$  is a terminal state.

(b) Reward func: ( $R$ )

$$R(s) = \begin{cases} -1 & \text{if } s \in \{s_0, s_1, s_2, s_3\} \\ 0 & \text{if } s = s_4. \end{cases}$$

Discount factor ( $\gamma$ )

$$\gamma = 1.$$

# Finding Expected no. of tories required.

Bellmann  
equ.

$$V(s) = E(r_{t+1}|s_t=s) + \gamma \sum_{s' \in S} P_{ss'} V(s')$$

$$V(s) = R(s) + \gamma \sum_{s' \in S} P_{ss'} V(s') \rightarrow (1)$$

# By calculating Eqn(1) for all states, we get

$$V(s_0) = -1 + (0.75 V(s_0) + 0.25 V(s_1))$$

$$V(s_1) = -1 + (0.75 V(s_0) + 0.25 V(s_1) + 0.25 V(s_2))$$

$$V(s_2) = -1 + (0.5 V(s_0) + 0.25 V(s_1) + 0.25 V(s_3))$$

$$V(S_3) = -1 + (0.5 V(S_0) + 0.25 V(S_1) + 0.25 V(S_4))$$

$$V(S_4) = \underline{0}$$

By solving the above eqns.

we get  $V(S_0) = -256$

$$V(S_1) = -252$$

$$V(S_2) = -240$$

$$V(S_3) = -192$$

$$V(S_4) = 0$$

# The value of  $R(s)$

# Since we choose  $R(s) = -1; \forall S \in \{S_0 \dots S_3\}$ , it represents that a step is counted whenever we move from one state other than at  $S_4$ , which is a terminal state. :-  $[V(S_0)]$  should determine the average no. of steps taken to reach the pattern.

$$\text{"Average" no. of steps} = 256/11$$

P2 Finite Horizon MDP

Reward :  $3s^2 + 5$

$s \in \{1, 2, 3, 4\}$  denoted as

$a \in \{\text{Continue, Quit}\} = \{C, Q\}$

$\gamma = 1$

@  $V_N(s) \dots ?$

Since at 'N', we need to quit anyways, we'll be left with the reward that we get

$$\therefore \boxed{V_N(s) = 3s^2 + 5}$$

$$V_N(1) = 8$$

$$V_N(2) = 17$$

$$V_N(3) = 32$$

$$V_N(4) = 53$$

(b)  $Q^{N-1}(s, a) \dots ?$

$$Q^{N-1}(s, Q) = R(s) = 3s^2 + 5 \rightarrow (1)$$

$$Q^{N-1}(s, C) = 0 + \underset{s}{E}(V^N(s)) = \frac{1}{4} \left( \sum_{s=1}^4 3s^2 + 5 \right) \rightarrow (2)$$

$$\Rightarrow Q^{N-1}(1, Q) = 8 \quad \left| \begin{array}{c} Q^{N-1}(1, C) \\ Q^{N-1}(2, C) \\ Q^{N-1}(3, C) \\ Q^{N-1}(4, C) \end{array} \right\} \frac{1}{4} (8 + 17 + 32 + 53) = 27.5$$

$$Q^{N-1}(2, Q) = 17$$

$$Q^{N-1}(3, Q) = 32$$

$$Q^{N-1}(4, Q) = 53$$

$$(c) \quad V^{N-1}(s) = \max_{a \in A} Q^{N-1}(s, a)$$

$$= \max(Q^{N-1}(s, Q), Q^{N-1}(s, C))$$

$$V^{N-1}(1) = \max(8, 27.5) \Rightarrow 27.5$$

$$V^{N-2}(2) = \max(17, 27.5) = 27.5$$

$$V^{N-2}(3) = \max(32, 27.5) = 32$$

$$V^{N-2}(4) = \max(53, 27.5) = 53$$

(d) For  $2 < n \leq N$

$$V^n(s) = \max_{a \in A} Q^n(s, a)$$

$$= \max(Q^{n-1}(s, Q), Q^{n-1}(s, C))$$

$$= \max(3s^2 + 5, E(V^n(s)))$$

(e) For  $2 < n \leq N$

$$Q^n(s, C) = \underset{\text{no reward}}{\underset{\text{if contd.}}{0}} + E_S(V^n(s))$$

$$= E_S(\max_{a \in A} Q^n(s, a))$$

$$= E_S(\max(Q^n(s, Q), Q^n(s, C)))$$

$$= E_S(\max(3s^2 + 5, Q^n(s, C)))$$

$$\Phi^n(s, c) = \underset{s'}{\mathbb{E}} (\max(3s^2 + 5, \hat{Q}^n(s', c')))$$

(f) Optimal policy:

$$\pi^n(s) = \begin{cases} Q: \text{quit} & ; \text{ if } \Phi^n(s, Q) \geq \Phi^n(s, C) \\ C: \text{continue} & ; \text{ o.w} \end{cases}$$

$\forall 2 \leq n \leq N$

(g) Open calculating  $\Phi^n(s, c)$ :  $\forall s \in \{1, 2, 3, 4\}$

$$\Phi^N(s, c) = 0$$

$$\Phi^{N-1}(s, c) = 27.5$$

$$\Phi^{N-2}(s, c) = 35$$

$$\Phi^{N-3}(s, c) = 39.5$$

$$\Phi^n(s, c) \in [35, 53] \quad \forall n \in [1, N-2]$$

$$\therefore \Phi^n(s, c) < 53$$

: The policy will be -

$$\forall n \in [1, N-2], \quad \pi_n^*(s) = c, \quad s \in \{1, 2, 3\}$$

$$\pi_n^*(s) = q; \quad s \in \{4\}$$

$$n = N-1; \quad \pi_n^*(s) = c; \quad s \in \{1, 2\}$$

$$\pi_n^*(s) = q; \quad s \in \{3, 4\}$$

$$n = N; \quad \pi(s) = q; \quad s \in \{1, 2, 3, 4\}$$

(g).  $\because \pi^*(s)$  depends upon ' $n$ '.

→ The optimal policy is non-stationary.

### Q3. Value Iteration

$$M \langle S, A, P, R, \gamma \rangle$$

$|S| < \infty$ ;  $|A| < \infty \rightarrow \gamma \in [0, 1)$

$$\hat{M} \langle S, A, P, \hat{R}, \gamma \rangle$$

$$|R(s, a, s') - \hat{R}(s, a, s')| \leq \epsilon$$

\* Given a policy  $\pi$ ,  $v^n, \hat{v}^n$  are policies

a) Expression relating  $v^n, \hat{v}^n$ ,

$$v^n(s) = \sum_a \pi(a|s) \cdot \sum_{s'} p_{ss'}^a (R_{ss'}^a + \gamma v^n(s'))$$

$$\hat{v}^n(s) = \sum_a \pi(a|s) \cdot \sum_{s'} p_{ss'}^a (\hat{R}_{ss'}^a + \gamma \hat{v}^n(s'))$$

Consider:

$$v^n(s) - \hat{v}^n(s) = \sum_a \pi(a|s) \sum_{s'} p_{ss'}^a (R_{ss'}^a - \hat{R}_{ss'}^a) \\ + r \sum_a \pi(a|s) \sum_{s'} p_{ss'}^a (v^n(s') - \hat{v}^n(s'))$$

absolute value on both sides  $\& |a+b| \leq |a| + |b|$   
on RHS

$$|v^n(s) - \hat{v}^n(s)| \leq \left| \sum_a \pi(a|s) \sum_{s'} p_{ss'}^a (R_{ss'}^a - \hat{R}_{ss'}^a) \right| \\ + \left| r \sum_a \pi(a|s) \sum_{s'} p_{ss'}^a (v^n(s') - \hat{v}^n(s')) \right|$$

$$\leq \sum \pi(q/s) \sum_{S'} P_{SS'}^q |R_{SS'}^q - \hat{R}_{SS'}^q|$$

$$+ r \sum_a \pi(a/s) \sum_S P_{SS'}^a |V_{(S')}^\pi - V_{(S)}^\pi|$$

(  $\because \pi$  &  $P$  are pos non-negative)

$$\leq \sum \pi(q/s) \sum_{S'} P_{SS'}^q (\varepsilon)$$

$$+ r \sum_a \pi(a/s) \cdot \sum_S P_{SS'}^a |V_{(S')}^\pi - V_{(S)}^\pi|$$

$$(\because \sum_a \pi(a/s) \sum_{S'} P_{SS'}^a (\varepsilon) = \sum_a \pi(a/s) \varepsilon = \varepsilon)$$

$$|V(s) - \hat{V}(s)| \leq \varepsilon + r \sum_a \pi(a/s) \sum_S P_{SS'}^a |V_{(S')}^\pi - V_{(S)}^\pi|$$

$$(b) V_x = \max_a \left( \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma \hat{V}_x(s')) \right)$$

$$\hat{V}_x = \max_a \left( \sum_{s' \in S} P_{ss'}^a (\hat{R}_{ss'}^a + \gamma \hat{V}_x(s')) \right)$$

$$V_x - \hat{V}_x = \max_a \left( \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V_x(s')) \right.$$

$$- \max_a \left( \sum_{s' \in S} P_{ss'}^a (\hat{R}_{ss'}^a + \gamma \hat{V}_x(s')) \right)$$

$$\max(a, b) \geq \max(a) - \max(b)$$

$$\leq \max_a \left( \sum_{s' \in S} P_{ss'}^a ((R_{ss'}^a - \hat{R}_{ss'}^a) + \gamma (V_x(s') - \hat{V}_x(s'))) \right)$$

$$\therefore R_{ss}^a - \hat{R}_{ss'}^a \leq \varepsilon$$

$$\leq \max_a \left( \sum_{s' \in S} P_{ss'}^a (\varepsilon) + \gamma (V_x(s') - \hat{V}_x(s')) \right)$$

$$V_x - \hat{V}_x \leq \max_a \left( \varepsilon + \gamma \sum_{s' \in S} P_{ss'}^a (V_x(s') - \hat{V}_x(s')) \right)$$

$\approx$

$$V_x - \hat{V}_x \leq \varepsilon + \max_a \left( \gamma \sum_{s' \in S} p_{ss'}^q (V_x^{(s')} - \hat{V}_x^{(s')}) \right)$$

(c) Yes, M & M' have the same optimal policy. Because,

$$V_x - \hat{V}_x \text{ are far by } \varepsilon \geq 0$$

$\therefore$  the  $\pi$  which achieves  $V_x$  also achieves  $\hat{V}_x$ .

$\Rightarrow$  thus both have same optimal policy

## P4. Effect of Noise & Discounting.

Need to map each given path to possible settings of  $\eta$  &  $\gamma$

$$\Rightarrow \eta = \{0, 0.5\}; \gamma \in \{0.1, 0.9\}$$

a) Close Exit (state with reward +1) but risk of cliff (dashed path to +1)

$$\rightarrow \boxed{\gamma = 0.1, \eta = 0}$$

Reason: with 'high' discounting factor discourages the agent to travel long paths as the final reward will be scaled down exponentially by  $\gamma$ .

→ Since ( $\eta = 0$ ): there is no randomness, hence, risk of falling into cliff is less.

(b) Distant exit (state with +10); risk the cliff (dashed path +10)

$$\rightarrow \boxed{\gamma = 0.9, \eta = 0}$$

Reason: with  $\gamma \neq 0$  → agent will prefer the distant route if the net reward it gets at

the end is high. Since higher  $\gamma$  will not significantly scale down the reward.

→  $\eta=0 \Rightarrow$  No noise  $\Rightarrow$  Reduces risk  
of falling into a cliff; can pick the path with cliff.

(c) Close exit (reward +1); avoid cliff  
(solid path +1)

$$\rightarrow \boxed{\gamma = 0.1 \quad \eta = 0.5}$$

$\gamma \downarrow \Rightarrow$  prefers shortest path although the reward at final state is low;

$\eta \uparrow \Rightarrow$  more randomness  $\Rightarrow$  more risky to fall into cliff  $\Rightarrow$  discourages encourages the agent to avoid cliff.

(d) Distant exit (reward +10); avoiding cliff  
(solid path +10)

$$\rightarrow \boxed{\gamma = 0.9 \quad \eta = 0.5}$$

→  $\gamma \uparrow \quad \eta \uparrow$ ; encourages to take distant path if net reward is greater.

→  $\eta \uparrow \Rightarrow$  risk of falling into cliff  $\Rightarrow$  encourages to avoid path with -ve reward.

P5. On Value Iteration Algo.

$M := \langle S, A, P, R, r \rangle ; |S| < \infty ; |A| \leq \infty , r \in [0, 1]$

→ Given: policy  $\pi$ ; need to evaluate  $V^\pi(s)$  using iterative update.

$$V_{k+1} \leftarrow \sum_a \pi(a|s) \sum_{s'} P_{ss'}^a [R_{ss'}^a + r V_k^{\pi}(s')] \\ \boxed{V_{k+1} = \frac{R + r P V_k^{\pi}}{(1 - r)}} \quad \text{(following notation from class)} \rightarrow (1)$$

while,

$$V_k^{\pi} \text{ Converge to } V^{\pi} = (I - rP)^{-1} R$$

$$\Rightarrow \boxed{V^{\pi} = R + r P V^{\pi}} \rightarrow (2)$$

⇒ Given, algo will terminate

$$\|V_{k+1} - V_k\|_{\infty} \leq \epsilon, \epsilon > 0. \rightarrow (3)$$

Q1) Need to prove:  $\|V_{k+1} - V^{\pi}\|_{\infty} \leq \frac{\epsilon r}{1-r}$

Consider LHS:

$$\|V_{k+1} - V^{\pi}\|_{\infty} = \|R + r P V_k - (R + r P V^{\pi})\|_{\infty}$$

$$= \|\gamma P(v_k - v^*)\|_\infty$$

$$\leq \gamma \|P\|_\infty \|v_k - v^*\|_\infty$$

$$\leq \gamma \|v_k - v^*\|_\infty$$

$$\leq \gamma \| (v_k - v_{k+1}) + (-v^* + v_{k+1}) \|_\infty$$

$$\leq \gamma (\|v_k - v_{k+1}\|_\infty + \|v_{k+1} - v^*\|_\infty)$$

( $\Delta k$  inequality)

$$\|v_{k+1} - v^*\|_\infty \leq \gamma (\|v_{k+1} - v_k\|_\infty + \|v_k - v^*\|_\infty)$$

$$(1-\gamma) \|v_{k+1} - v^*\|_\infty \leq \gamma \|v_{k+1} - v_k\|_\infty$$

$$\leq \gamma \varepsilon \quad (\because \text{From Eq (3)})$$

$$\boxed{\|v_{k+1} - v^*\|_\infty \leq \frac{\gamma \varepsilon}{1-\gamma}}$$

~~(b)~~ (b) Need to prove

$$\|v_{k+1} - v^*\|_\infty \leq \gamma^k \|v_1 - v_\infty^*\|.$$

Substituting using Eq (1) & (2)

Consider LHS:

$$\begin{aligned}\|v_{k+1} - v^*\|_\infty &= \|(R + \gamma P V_k) - (R P V^*)^\pi\|_\infty \\ &= \|\gamma P(v_k - v^*)\|_\infty \\ &\leq \gamma \|P\|_\infty \|v_k - v^*\|_\infty \\ &\leq \gamma \|v_k - v^*\|_\infty \quad (\text{proved in class})\end{aligned}$$

$$\therefore \|v_{k+1} - v^*\|_\infty \leq \gamma \|v_k - v^*\|_\infty \rightarrow (a)$$

Similarly

$$\begin{aligned}\Rightarrow \|v_k - v^*\|_\infty &\leq \gamma \|v_{k-1} - v^*\|_\infty \\ \|v_{k-1} - v^*\|_\infty &\leq \gamma \|v_{k-2} - v^*\|_\infty \\ &\vdots \\ \|v_2 - v^*\|_\infty &\leq \gamma \|v_1 - v^*\|_\infty\end{aligned}\quad \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} k-1 \text{ inequalities}$$

$$\Rightarrow \|v_k - v^*\|_\infty \leq \gamma^{k-1} \|v_1 - v^*\|_\infty$$

Substituting in eq (a)

$$\begin{aligned}\|v_{k+1} - v^*\|_\infty &\leq \gamma (\gamma^{k-1} \|v_1 - v^*\|_\infty) \\ &\leq \gamma^k \|v_1 - v^*\|_\infty\end{aligned}$$

Hence proved

(c)  $v$ : denote the value func.  
acts 's' → define the states

$S \in \mathbb{S}$ :  $\rightarrow \mathbb{S}$ : state space  
 $|S| = d$  : d states

then  $v$  is represented by a vector of d-dimension

$$v: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$C: L(v) = \max_{a \in A} [R^a + \gamma P^a v]$$

$$L: V \rightarrow V : L \text{ also maps } \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$\Rightarrow$  Need to prove:  
 $L(u) \leq L(v)$ ; if  $u \leq v$

$$u \leq v \Rightarrow u_i \leq v_i : i \in \{1, \dots, d\}$$

$\Rightarrow$  Consider; the following:

$$v \geq u \quad (\text{as defined})$$

Consider an arbitrary action "a"  $\in A$ ; st.

$$P^a_v \geq P^a_u \quad (\text{holds true since } P^a \text{ contains elements ranging from } [0,1])$$

$$\gamma P^a_v > \gamma P^a_u \quad (\because \gamma \in [0,1])$$

$$R^a + \gamma P^a v \geq R^a + \gamma P^a u$$

(adding  $R^a$  for both sides)

$\Leftrightarrow (1)$

$\Rightarrow$  Now consider the following

$$\max_{a \in A} (R^a + \gamma P^a v) \geq R^a + \gamma P^a v \quad (\text{defn of maximum})$$

$\hookrightarrow$  at  $A$ .

$$\geq R^a + \gamma P^a u \quad (\text{From eq. (1)})$$

$$\geq \max_{a \in A} (R^a + \gamma P^a u) \quad (\because \text{By defn of max operator})$$

$$\therefore \max_{a \in A} (R^a + \gamma P^a u) \geq \max_{a \in A} (R^a + \gamma P^a v)$$

$$L(v) \geq L(u)$$

$$\Rightarrow L(w) \leq L(v) \text{ if } u \leq v$$

Hence, Bellman optimality operator is

monotone

Pb. On Contractions

@ P, Q contractions on normed vector space  $\langle V, \|\cdot\| \rangle$ .

i.e;  $P: V \rightarrow V$ ,  $Q: V \rightarrow V$

$\forall u, v \in V$ ;  $r_p, r_q$  ;  $r \in [0, 1)$

$$\|P(v) - P(u)\| \leq r_p \|v - u\| \rightarrow (1)$$

$$\|Q(v) - Q(u)\| \leq r_q \|v - u\| \rightarrow (2)$$

By the defn of contraction

Prooving  $P \circ Q$ ,  $Q \circ P$  are contractions on the same vector space.

Consider vector space  $\langle V, \|\cdot\| \rangle$ .

$P \circ Q : V \rightarrow V$  (since, range of Q is V and domain of P is V)

$\therefore P \circ Q : V \rightarrow V$

$\forall u, v \in V$

$$\|P(Q(v)) - P(Q(u))\| \leq r_p \|Q(v) - Q(u)\|$$

$$\leq r_p \cdot r_q \|v - u\|$$

$$= r_p r_q \|v - u\|$$

→ using eq(2)

$$\therefore \|P(Q(v)) - P(Q(u))\| \leq r_p r_q \|v - u\|$$

Since  $r_p, r_q \in [0, 1]$ ,

$$\eta = r_p r_q \in [0, 1]$$

$P \circ Q$  is a contraction, with  $\eta_{pq}$  as Lipsitz coefficient, &  $\eta_{pq} \in [0, 1]$

# Consider  $Q \circ P$

$\boxed{Q \circ P : V \rightarrow V}$  ( $\because$  Range of  $P$  is  $V$  & Domain of  $Q$  is  $V$ )  
Range of  $Q$  is  $V$

consider  $v, u \in V$

$$\begin{aligned} \|Q(P(v)) - Q(P(u))\| &\leq r_q \|P(v) - P(u)\| \\ &\leq r_q (r_p \|v - u\|) \\ &= r_p r_q \|v - u\| \quad (\text{From eq (1)}) \end{aligned}$$

$$\|Q(P(v)) - Q(P(u))\| \leq r_q r_p \|v - u\|$$

$$\because r_p, r_q \in [0, 1] \Rightarrow$$

$$\Rightarrow \|Q(P(v)) - Q(P(u))\| \leq \eta_{qp} \|v - u\|; \quad \eta_{qp} \in [0, 1]$$

Hence,  $Q \circ P$  is also a contraction on the

same normal vector space  $V$ . with coefficient

$$\eta_{qp} \in [0, 1]$$

(b) As calculated earlier,

Lipschitz coefficient for  $P \circ Q = Q \circ P = \underline{\underline{\gamma_p \gamma_q}}$

where  $\gamma_p$  is the coeff for  $P$ :

$\gamma_q$  . . . for  $Q$ .

(c) given  $B := F \circ L$

$\Rightarrow L$  is the Bellman optimality operator -

It has been proved that,  $L$  is a contraction.  
(in class).

$\therefore L$  is a contraction, and is defined on a normal vector space denoted by  $\langle \mathcal{V}, \|\cdot\| \rangle$

$$L : \mathcal{V} \rightarrow \mathcal{V}$$

2) Conditions for  $B$  to converge.

a)  $F$  should be defined on the same normal vector space as  $L$  i.e;  $\underline{\mathcal{V}}$

b)  $F$  should be a contraction.

c) The initial guess in the value iteration algo  
(lets say)  $v \in \underline{\mathcal{V}}$

Reason:

→ Banach Fixed point theorem states that;

& if for normal vector space  $\langle \mathcal{V}, \|\cdot\| \rangle$ , a contraction  $B : \mathcal{V} \rightarrow \mathcal{V}$ ; Then the iterative

application of  $B$  converges to a fixed point in  $\mathcal{V}$  (thus resulting in a unique soln).

$\therefore$  For  $B: F_0 L$  should be a contraction  
on the same vector space of  $L$  i.e.,  $\underline{V}$

$\rightarrow$  For that to happen, we proved in (6a)  
that if  $F, L$  defined on the same <sup>normed</sup> vector space  
 $\langle V, \|\cdot\| \rangle$  and are contractions,  
then  $B = F \circ L$  is also defined on  $\langle V, \|\cdot\| \rangle$   
and is also a contraction.  
 $\Rightarrow$  Thus  $B$  satisfies Banach Fixed Point thm.  
& thus converges when the above conditions  
are satisfied

