

Avinash. G. Kon
ED15B006.

$$\textcircled{1} \quad f: \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad f \in C_1^{0,0}(\Omega)$$

$$\Rightarrow \|f(x) - f(y)\| \leq L\|x - y\|. \rightarrow \text{(i)}$$

A10. $y = x^*$ $L = \alpha$.

$$\Rightarrow \|f(x) - f(x^*)\| \leq \alpha\|x - x^*\|.$$

$$\Rightarrow \|f(x) - x^* + x - x^*\| \leq \alpha\|x - x^*\|$$

$$\Rightarrow \|f(x) - x\| + \|x - x^*\| \leq \alpha\|x - x^*\|. \quad (\because \|a+b\| \leq \|a\| + \|b\|)$$

$$\Rightarrow \|f(x) - x\| \leq (\alpha-1)\|x - x^*\|.$$

$$\Rightarrow \|x - x^*\| \geq \frac{1}{(\alpha-1)}\|f(x) - x\|.$$

$$\Rightarrow \left[\|x^* - x\| \leq \frac{1}{(1-\alpha)}\|f(x) - x\| \right] \rightarrow \underline{\text{Proved}}$$

b) f is monotone $\Rightarrow f(x) > f(y) \quad \forall x > y$.

given $x \leq y + \delta e$

$$\Rightarrow f(x) \leq f(y + \delta e)$$

$\Rightarrow f \in C_1^{0,0}(\Omega)$

$$\Rightarrow \|f(y + \delta e) - f(y)\| \leq \alpha\|\delta e\|.$$

$$\Rightarrow -\alpha\delta e \leq f(y + \delta e) - f(y) \leq \alpha\delta e$$

$$\Rightarrow f(y) - \alpha\delta e \leq f(y + \delta e) \leq f(y) + \alpha\delta e$$

$$\Rightarrow f(y + \delta e) \leq \alpha\|\delta e\| + f(y)$$

$$\Rightarrow \boxed{f(x) \leq \alpha\|\delta e\| + f(y)} \rightarrow \text{proved.}$$

$$\textcircled{2} \quad X = \{A, B\} \quad a = \{\text{cont}, \text{toggle}\}$$

$$R = \{r_A, r_B; c\}; P_{ij}(\text{cont}) = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}; P_{ij}(\text{toggle}) = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

Let π_1 be stationary policy $\pi_1 = \{\pi_{1A}, \pi_{1B}\}$

- for $\pi_1' = \{\text{cont}, \text{cont}\}$

i.e., to continue to work in state they are in...

$$\Rightarrow J_{\pi_1'}(A) = r_A + \alpha J_{\pi_1'}(A)$$

$$\boxed{J_{\pi_1} = J_{\pi_1'}, J_{\pi_1}}$$

$$\Rightarrow J_{\pi_1'}(A) = \frac{r_A}{1-\alpha}$$

$$\text{IIIY } J_{\pi_1'}(B) = r_B + \alpha J_{\pi_1'}(B)$$

$$\Rightarrow J_{\pi_1'}(B) = \frac{r_B}{1-\alpha}$$

$$\overline{J_{\pi_1'} = \left[\frac{r_A}{1-\alpha}, \frac{r_B}{1-\alpha} \right]} \rightarrow \text{Policy evaluation for } \pi_1'$$

- for $\pi_1^2 = \{\text{cont}, \text{toggle}\}$

$$J_{\pi_1^2}(A) = r_A + \alpha J_{\pi_1^2}(A) \Rightarrow \overline{J_{\pi_1^2}(A) = \frac{r_A}{1-\alpha}}$$

$$J_{\pi_1^2}(B) = r_A - C + \alpha J_{\pi_1^2}(A) \Rightarrow \overline{J_{\pi_1^2}(B) = \frac{r_A}{1-\alpha} - C}$$

$$\overline{J_{\pi_1^2} = \left[\frac{r_A}{1-\alpha}, \frac{r_A}{1-\alpha} - C \right]} \rightarrow \text{Policy evaluation for } \pi_1^2$$

- for $\pi^3 = \{\text{toggle, work}\}$,
 i.e. jump to state B if in state A,
 and continue to stay in state B if in state B.

$$J_{\pi,3}(A) = r_B - c + \alpha J_{\pi,3}(B)$$

$$J_{\pi,3}(B) = r_B + \alpha J_{\pi,3}(B) \Rightarrow J_{\pi,3}(B) = \frac{r_B}{1-\alpha}$$

$$\Rightarrow J_{\pi,3}(A) = r_B - c + \frac{\alpha r_B}{1-\alpha} = \frac{r_B}{1-\alpha} - c.$$

$$\underline{J_{\pi,3} = \left[\frac{r_B}{1-\alpha} - c, \frac{r_B}{1-\alpha} \right]} \rightarrow \begin{array}{l} \text{Policy evaluation} \\ \text{for } \pi^3 \end{array}$$

- for $\pi^4 = \{\text{toggle, toggle}\}$,
 i.e. change state irrespective of current state.

$$J_{\pi,4}(A) = r_B - c + \alpha J_{\pi,4}(B)$$

$$J_{\pi,4}(B) = r_A - c + \alpha J_{\pi,4}(A)$$

$$\Rightarrow J_{\pi,4}(A) = \frac{r_A + \alpha r_B - (1+\alpha)c}{1-\alpha^2}$$

$$J_{\pi,4}(B) = \frac{r_B + \alpha r_A - (1+\alpha)c}{1-\alpha^2}$$

$$\Rightarrow \underline{J_{\pi,4} = \left[\frac{r_A + \alpha r_B - (1+\alpha)c}{1-\alpha^2}, \frac{r_B + \alpha r_A - (1+\alpha)c}{1-\alpha^2} \right]} \rightarrow \begin{array}{l} \text{Policy} \\ \text{evaluation} \\ \text{for } \pi^4 \end{array}$$

$\hat{y} \alpha \rightarrow 0$

$J_{T_1^1} = [r_A, r_B]$ will be optimal

Since in other policy cost at few states will be $\underline{-ve}$.

$\Rightarrow \alpha \rightarrow 0 \quad \tau_1^1 = d(\text{cont}), \text{ write } \}$ is the best policy.

As $\alpha \rightarrow 1$.

$$J_{T_1^1} \cdot (1-\alpha) = [r_A, r_B] ; J_{T_1^2} \cdot (1-\alpha) = [r_A, r_A]$$

$$J_{T_1^3} \cdot (1-\alpha) = [r_B, r_B] ; J_{T_1^u} \cdot (1-\alpha) = \left[\frac{r_A+r_B-c}{2}, \frac{r_A+r_B-c}{2} \right]$$

As compared to all the J 's

$J_{T_1^2}$ is optimal wst. since $\underline{r_A > r_B}$

$\Rightarrow \tau_1^2 = d(\text{cont}), -\text{logit} \}$ move to state A

is best policy for $\alpha \rightarrow 1$.

b) $c = 3 \quad r_A = 2 \quad r_B = 1 \quad \alpha = 0.9$.

$$J_{T_1^1} = [20, 10] ; J_{T_1^2} = [20, 17]$$

$$J_{T_1^3} = [7, 10] ; J_{T_1^u} = [-15.26, -14.74]$$

from the above values it clear that moving to state A is optimal.

(3)

$$m_j^* = \min_i \min_a P_{ij}(a)$$

$$\tilde{P}_{ij} = \frac{P_{ij} - m_j^*}{1 - \sum_k m_{ik}}$$

\Rightarrow for \tilde{P}_{ij} to be transition probabilities

- $0 < \tilde{P}_{ij} < 1$. $\sum_j \tilde{P}_{ij} = 1$ $\forall i$

$$\sum_j \tilde{P}_{ij} = \sum_j \frac{P_{ij} - m_j^*}{1 - \sum_k m_{ik}} = \frac{1}{(1 - \sum_k m_{ik})} \sum_j (P_{ij} - m_j^*)$$

$$\Rightarrow \frac{\sum_j (P_{ij} - m_j^*)}{1 - \sum_k m_{ik}} = \frac{1 - \sum_k m_{ik}}{1 - \sum_k m_{ik}} = 1. \quad \left\{ \begin{array}{l} \text{• } P_{ij} \text{ is} \\ \text{transition} \\ \text{probability} \end{array} \right.$$

$$m_j^* \leq P_{ij} \quad \left\langle \begin{array}{l} \text{• } m_j^* \text{ involves min. over } P_{ij} \end{array} \right\rangle$$

$$\Rightarrow \frac{P_{ij} - m_j^*}{1 - \sum_k m_{ik}} \geq 0 \quad \left\langle \begin{array}{l} \text{• } 1 - \sum_k m_{ik} > 0 \text{ given} \end{array} \right\rangle$$

We know that:

$$P_{ij} - m_j^* \leq \sum_j (P_{ij} - m_j)$$

$$\Rightarrow \frac{P_{ij} - m_j^*}{1 - \sum_k m_{ik}} \leq 1. \quad \rightarrow \text{(ii)}$$

By (i) & (ii)

$$0 < \tilde{P}_{ij} \leq 1$$

proved

$$b) J^*(x) = \lim_{N \rightarrow \infty} T^N J(x) \quad \forall x \in S \quad (6)$$

W.S.B.C. set of all states

$|g(i, a)| \leq M \Leftrightarrow g(i, a)$ is bounded.

$$\overline{J}^*(x) = \min_a \lim_{N \rightarrow \infty} \sum_{k=0}^{N-1} (\alpha \sum_{j} P_{ij}^k(a))^k g(x_k, a)$$

$$\overline{J}^*(x) = \min_a [g(i, u) + \alpha \sum_j P_{ij}(a) \overline{J}^*(x_j)] \quad \forall x_i, x_j \in S. \quad (1)$$

$$\begin{aligned} \text{and } \widetilde{J}(x_i) &= \min_a [g(i, u) + \alpha \sum_j \tilde{P}_{ij}(a) \widetilde{J}(x_j)] \\ &= \min_a [g(i, u) + \alpha (1 - \sum_k m_k) \sum_j \frac{(P_{ij} - m_j) \widetilde{J}(x_j)}{1 - \sum_k m_k}] \\ &= \min_a [g(i, u) + \alpha \sum_j (P_{ij} - m_j) \widetilde{J}(x_j)] \\ &= \min_a [g(i, u) + \alpha \sum_j P_{ij}(a) \widetilde{J}(x_j) + \alpha \sum_j P_{ij}(a) \\ &\quad - \alpha \sum_j m_j \widetilde{J}(x_j)] \end{aligned}$$

adding $\alpha \frac{\sum_k m_k \widetilde{J}(x_k)}{1 - \alpha}$ on both sides

$$\Rightarrow \widetilde{J}(x_i) + \alpha \frac{\sum_k m_k \widetilde{J}(x_k)}{1 - \alpha} = \min_a [g(i, u) + \alpha \sum_j P_{ij}(a) \widetilde{J}(x_j) - \left(\alpha - \frac{\alpha}{1-\alpha}\right) \sum_k \widetilde{J}(x_k) m_k]$$

$$\begin{aligned} &= \min_a [g(i, u) + \alpha \sum_j P_{ij}(a) \widetilde{J}(x_i) + \underline{\alpha^2} \sum_j P_{ij}(a) \sum_k \widetilde{J}(x_k) m_k] \end{aligned}$$

$\left[\because \sum_j P_{ij} = 1 \right] \quad \leftarrow$

$$\Rightarrow \tilde{J}(x_i) + \alpha \frac{\sum_k m_k \tilde{J}(x_k)}{1-\alpha} = \min_{\sigma} [g(i, \sigma) + \alpha \sum_j P_{ij}^{\sigma}(0) \left[\tilde{J}(x_i) + \alpha \frac{\sum_k m_k \tilde{J}(x_k)}{1-\alpha} \right]]$$

By comparing ② with ①

\longrightarrow ②

$$J^*(x_i) = \tilde{J}(x_i) + \alpha \frac{\sum_k m_k \tilde{J}(x_k)}{1-\alpha}$$

$$\Rightarrow J^* = \tilde{J} + \alpha \frac{\sum_k m_k \tilde{J}(k)}{1-\alpha} \longrightarrow \underline{\text{proved.}}$$

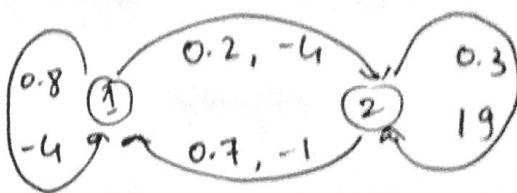
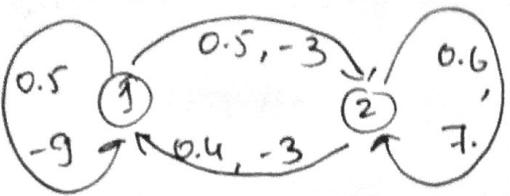
$$\text{As } J^* = \tilde{J} + C$$

$$\text{where } C = \alpha \frac{\sum_k m_k \tilde{J}(x_k)}{1-\alpha}$$

Policy minimizing J^* will
doubtly minimize \tilde{J}

\Rightarrow \tilde{J} & \tilde{J} has some optimal policy.

$$X = \{1, 2\} ; A = \{a, b\} ; \alpha = 0.9.$$



action a.

$$\underline{a \neq b} \quad \pi_0 = \{a, a\}$$

Policy evaluation @ π_0

$$\begin{aligned} J_{\pi_0}(1) &= 0.5(-9 + \alpha J_{\pi_0}(1)) + 0.5(-3 + \alpha J_{\pi_0}(2)) \\ &= -6 + 0.5\alpha [J_{\pi_0}(1) + J_{\pi_0}(2)] \end{aligned}$$

$$J_{\pi_0}(2) = 3 + 0.4\alpha J_{\pi_0}(1) + 0.6\alpha J_{\pi_0}(2)$$

$$J_{\pi_0} = [-15.5, -5.6]$$

Policy update ($\pi_0, J_{\pi_0} = T J_{\pi_0}$)

$$\begin{aligned} T J_{\pi_0}(1) &= \min \{-6 + 0.5\alpha[-21.1], -4 + \alpha(0.8(-15.5) \\ &\quad + 0.2(-5.6))\} \\ &= \min \{-15.56, -16.2\} \\ &= -16.2 \Rightarrow \pi_1(1) = b. \end{aligned}$$

$$\begin{aligned} T J_{\pi_0}(2) &= \min \{3 + \alpha(0.4(-15.5) + 0.6(-5.6)), \\ &\quad 5 + \alpha(0.3(-5.6) + 0.7(-15.5))\} \\ &= \min \{-5.56, -6.313\} \\ &= -6.313 \Rightarrow \pi_1(2) = \underline{b} \end{aligned}$$

$$\pi_1 = \underline{\{b, b\}}$$

Policy evaluation @ d b, b

$$J_{\pi_1}(1) = 0.8(-4) + 0.2(-6) + \alpha [0.8 J_{\pi_1}(1) + 0.2 J_{\pi_1}(2)]$$

$$J_{\pi_1}(2) = 0.7(-1) + 0.3(19) + \alpha (0.7 J_{\pi_1}(1) + 0.3 J_{\pi_1}(2))$$

$$\Rightarrow J_{\pi_1} = \{-22.2, -12.3\}$$

Second update ($\pi_{\pi_1}, J_{\pi_1} = \pi J_{\pi_1}$)

$$\begin{aligned}\pi J_{\pi_1}(1) &= \min \{-6 + 0.5\alpha[-34.5], -4 + \alpha[0.8(-22.2) + 0.2(-12.3)]\} \\ &= \min \{-21.52, -22.24\} \\ &= \underline{-22.2} ; \quad \underline{\pi_1(1) = b}\end{aligned}$$

$$\begin{aligned}\pi J_{\pi_1}(2) &= \min \{3 + \alpha(0.4(-22.2) + 0.6(-12.3)), \\ &\quad 5 + \alpha(0.3(-12.3) + 0.7(-22.2))\} \\ &= \min \{-11.6, -12.7\} \\ &= \underline{-12.7} ; \quad \underline{\pi_1(2) = b}\end{aligned}$$

$$\Rightarrow \pi_1 = \underline{d b, b}$$

$\therefore \pi_1 = \underline{\pi_1} = \underline{d b, b}$ π_1 is an optimal policy

(6)

$$\underline{Q} \quad J = (0, 0)$$

$$TJ = \min_a \sum P_{ij}(a) [g(i, a, j) + \alpha J(j)]$$

$$TJ(1) = \min [-6, -4] = -6.$$

$$TJ(2) = \min \{3, 5\} = 3.$$

$$\rightarrow TJ = \underline{[-6, 3]} \Rightarrow \underline{\pi_i(a, a)}$$

$$T^2J = \min_a \sum P_{ij}(a) [g(i, a, j) + \alpha TJ(j)]$$

$$T^2J(1) = \min [-6 + 0.4(-3), -4 + 0.8(0.8(-6) + 0.2(3))] \\ = \min [-7.35, -7.78] = \underline{-7.78} \quad (b)$$

$$T^2J(2) = \min [3 + \alpha(0.4(-6) + 0.6(3)), 5 + \alpha(0.3(3) + 0.7(-6))] \\ = \min [2.46, 203] \Rightarrow \underline{203} \quad (b).$$

$$T^2J = \underline{[-7.78, 203]} \quad \pi_i = \underline{\{b, b\}}$$

$$\rightarrow T^3J(1) = \min [-6 + 0.41(-5.78), -4 + 0.8(0.8(-7.78) + 0.2(203))] \\ = \min [-8.58, -9.23] = \underline{-9.23}$$

$$T^3J(2) = \min [3 + \alpha(0.4(-7.78) + 0.6(203)), \\ 5 + \alpha(0.3(-9.23) + 0.7(-8.58))] \\ = \min [1.3, 0.06] = \underline{0.06}$$

$$T^3J = \underline{[-9.23, 0.06]} \quad \pi_i = \underline{\{b, b\}}$$

$$\begin{aligned} \pi^4 J(2) &= \min \{-6 + 0.4(-9.23), -4 + 0.4(0.8(-9.23) + 0.2(0.06))\} \\ &= \min [-10.12, -10.638] = \underline{-10.638} \\ \pi^4 J(2) &= \min [3 + 0.4(-9.23) + 0.6(0.06)], \\ &\quad 5 + 0.4[0.7(-9.23) + 0.3(0.06)] \end{aligned}$$

$$= \min [-0.29, -0.80] = \underline{-0.8}$$

$$\pi^4 J = [-10.63, -0.80] \quad \pi_1 = \underline{\{b, b\}}$$

iii. finding optimal policy with $\alpha = \underline{0.1}$

$$\text{let } \pi_{1,0} = \{a, a\}$$

$$J_{\pi_{1,0}}(1) = -6 + 0.05(J_{\pi_{1,0}}(1) + J_{\pi_{1,0}}(2))$$

$$J_{\pi_{1,0}}(2) = 3 + 0.04 J_{\pi_{1,0}}(1) + 0.06 J_{\pi_{1,0}}(2)$$

$$J_{\pi_{1,0}} = [-6.16, 2.92]$$

Policy update

$$\begin{aligned} \pi J_{\pi_{1,0}}(1) &= \min \{-6 + 0.05(-3.24), -4 + 0.05(0.8(-6.16) + 0.2(2.92))\} \\ &= \min \{-6.16, 2.92\} = \underline{-6.16} \end{aligned}$$

$$\begin{aligned} \pi J_{\pi_{1,0}}(2) &= \min \{3 + 0.1(0.4(-6.16) + 0.6(2.92)), \\ &\quad 5 + 0.1(0.3(+2.92) + 0.7(-6.16))\} \\ &= \min \{2.9, 4.6\} = 2.9 \end{aligned}$$

$$\Rightarrow \pi_1 = \{a, a\} = \pi_{1,0} \Rightarrow \{a, a\} \text{ is optimal.}$$

As $\alpha \rightarrow 0 \Rightarrow$ policy biased towards immediate reward / cost.

- @ state ① immediate minimum cost is to select action ② & stay in state ①.
- @ state ② immediate minimum cost is achieved by selecting action ① & jump to state ①