

DECLARATION.

I pledge that I have not copied or given any unauthorised assistance or this assignment.

ACKNOWLEDGEMENT :

- Bertsekas DPOC vol 1 Chap 1
- Class notes (notes taken by me during class & notes not used for teaching)
- Classmate :
 - i) Aniswar & (CS18B080), A~~T~~S Abhishek (FE18B001)
 - discussed Q3
 - ii) Narendhiran (CH18B015)
 - ^{wrong} checked Q4 transition probability ,
but gave hints for Q5

①
a)

Action space : { forage patch 1, forage patch 2,
 (call them as a_1, a_2, a_3 respectively) forage patch 3 }

State Space : { animal 2 units of 3 units of
 is dead, energy reserve, energy in
 reserve, 4 units of energy }
 in reserve }

We want to maximise probability of survival at
 the end of 3 time periods.

$$g_3(x_3) = \begin{cases} 1 & \text{if } x_3 \in \{2, 3, 4\} \\ 0 & \text{if } x_3 = \text{dead.} \end{cases}$$

Terminal cost

So we give a terminal cost of 1 if the
 animal has survived & all other transitions
 costs to be 0. (as we want to maximise P)

$$\text{i.e. } g_k(x_i, a_i, j) = 0 + \frac{1}{k+1} \delta_{ij}$$

So by maximising J_1

$$J = \max_k E(g_k(x_0))$$

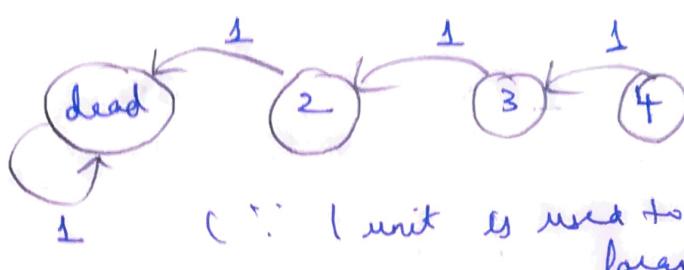
$$= \max_k \Theta(g_3(x_3))$$

we in fact maximise the probability of survival
 $(= E(g_3(x_0)))$

$$\sum_{k=0}^{n-1} g_k(x_k, u_k(x_k), x_{k+1})$$

Let us note the transition probabilities for each action

a1: Forage patch 1



(\because 1 unit is used to forage)

$$\begin{aligned} P_{dd} &= 1 \\ P_{2d} &= 1 \quad P_{ij} = 0 \\ P_{32} &= 1 \quad \text{for all} \\ P_{43} &= 1 \quad \text{others} \end{aligned}$$

a2: Forage patch 2

Assume: Risk of predation independent of finding food.

\therefore risk of predation = 0.004, there is a 0.996 probability that the animal survives

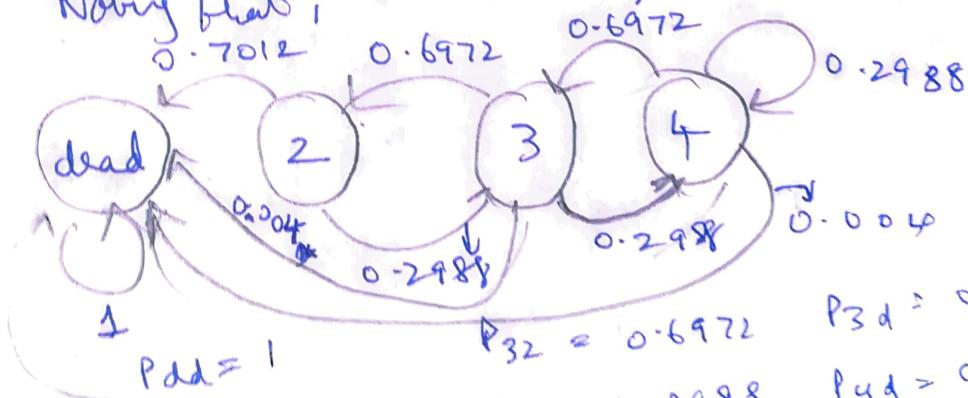
In the 0.996 chance the animal survives, there

is 0.3 chance to get food of 2 energy units
 $P(\text{food} | \text{surv}) = 0.3 \Rightarrow P(\text{food} \& \text{surv}) = 0.3 \times P_{\text{surv}}$
 (gaining 1 energy unit as a result - 1 unit will be used to forage)

& 0.7 chance it doesn't get food

(\Rightarrow loses 1 unit as a result of foraging)

$$\text{Noting that, } 0.996 \times 0.3 = 0.2988 \text{ & } 0.996 \times 0.7 = 0.6972$$



$$P_{dd} = 1$$

$$P_{2d} = 0.7012$$

$$P_{23} = 0.6972$$

$$P_{32} = 0.6972$$

$$P_{34} = 0.2988$$

$$P_{43} = 0.6972$$

$$P_{3d} = 0.004$$

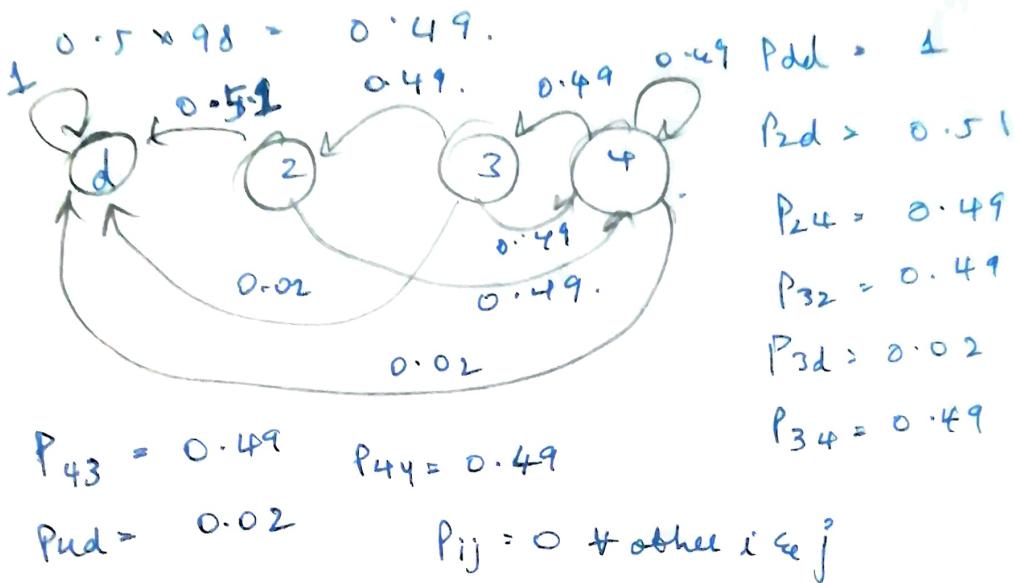
$$P_{4d} = 0.004$$

$$P_{ij} = 0 \text{ for all } i \neq j$$

as: Foreign patch 3

Here the busy gain 3 units (\Rightarrow net gain of 2 units)

Procedure similar to the case of α_2



$$\text{b) } J_3^*(x_3) = g_3(x_3) \\ = \begin{cases} 1 & \text{if } x_3 \in \{2, 3, 4\} \\ 0 & \text{if } x_3 = d. \end{cases}$$

DP Algorithm:

$$J_2^*(x_2) = \max_a E_{x_3} \left(g_2(x_2, a, x_3) + J_3^*(x_3) \right) \\ = \max_a \sum_{j \in S} P_{ij}(a) (g_2(x_2, a, j) + J_3^*(j))$$

but $g(1, a, j) > 0$ always

$$\Rightarrow J_2^*(x_2) = \min_a \sum_j P_{ij}(a) J_3^*(j)$$

~~but~~ $x_2 = d \Rightarrow$ action doesn't matter the animal
is already dead

action 1: forage patch 1

$$x_2 \quad \sum J^+ P_i j$$

d	0			
---	---	--	--	--

$$2 \quad J + 0 = 0$$

$$3 \quad P_{32} J_3^+(2) = 1$$

$$4 \quad P_{43} J_3^+(3) = 1$$

action 2: forage patch 2

$$x_2 \quad \sum J^+ P_i j$$

d	0			
---	---	--	--	--

$$2 \quad 0.2988 J_3^+(3) = 0.2988$$

$$1 \quad 0.2988 J_3^+(4) = 0.2988 \cdot 0.996$$

$$3 \quad + 0.6972 J_3^+(2)$$

$$4 \quad 0.2988 J_3^+(4) = 0.996$$

$$+ 0.6972 J_3^+(2)$$

action 3: forage patch 3

$$x_2 \quad \sum J^+ P_i j$$

d	0			
---	---	--	--	--

$$d \quad 0$$

$$x_2 \sum_{j=1}^3 P_{2j} J_2^+(j)$$

$$2 \quad P_{24} J_2^+(4) = 0.49.$$

$$3 \quad 0.49 J_2^+(2) + 0.49 J_2^+(4) = 0.98$$

$$4 \quad 0.49 J_2^+(3) + 0.49 J_2^+(4) = 0.98$$

choosing the actions that minimise the probability next $J_2(x_2)$ we have

x_2	$\mathbb{E}_2 P J_2^+(x_2)$	a. — (down + netty)
d	0	
2	0.49.	a_3 (foreign patch 3)
3	1	a_1 (foreign patch 1)
4	1	a_1 (foreign patch 1)

$$\text{Now, } J_1^+(x_1) = \max_a \mathbb{E} (g(x_1, a, j) + J_2^+(j)) \\ = \max_a \mathbb{E} (J_2^+(j|a))$$

If we chose action a ,

$$\mathbb{E}(J_2^+(j)) = \sum p_{ij}(a) J_2^+(j)$$

We make a table & compute $J_1(x_1)$ values for each action and then decide the optimal value

x_1	a_1	a_2	a_3
d	0	0	0
2	0	$\frac{0.2988 J_2^+(3)}{= 0.2988}$	$\frac{0.49 J_2^+(4)}{= 0.49}$
3	$1 \times J_2^+(2)$ $= 0.49$	$0.2988 J_2^+(4)$ $+ 0.6472 J_2^+(2)$ $= 0.6404$	$0.49 J_2^+(4)$ $+ 0.49 J_2^+(3)$ $= 0.7301$
4	$1 \times J_2^+(3)$ $= 1$	0.996	$0.49 J_2^+(4)$ $+ 0.49 J_2^+(3)$ $= 0.98$

Choosing the action that maximizes probability of return,

x_1	$J_1^+(x_1)$	Action
d	0	- (can't take any action - dead)
2	0.49	Forage patch 3 (a_3)
3	0.7301	Forage patch 3 (a_3)
4	1	Forage patch 1 (a_1)

$$\text{Similarly, } J_0^*(x_0) = \max_a \bar{E}(J_1^*(x_1)) \\ = \max_a \sum_{Hj} p_{ij}(a) (f_i^*(x_1))$$

We once again make a table of action & reward

x_0	$J_0^+(x_0)$	a_2	a_3
d	0	0	0
2	0	$0.2998(0.73)$ $= 0.218$	$0.49 + \cancel{0.49}$
3	0.49	$0.2998(1)$ $+ 0.6972(0.49)$ $= 0.6404$	$0.49 + (0.49)$ $= 0.7301$
4	0.7301	$(0.2998) + (0.6972)(0.73)$ $= 0.8078$	$0.49 +$ $0.49 \cancel{+ 0.7301}$ $= 0.8477$

Choosing the action that optimises the reward,

x_0	$J_0^+(x_0)$	Action
d	0	For - (dead)
2	0.49	Forage patch 3 (a_3)
3	0.7301	Forage patch 3 (a_3)
4	0.8477	Forage patch 3 (a_3)

Optimal policy for foraging obtained from the analysis!

$$\pi^* = \{ \pi_0(x), \pi_1(x), \pi_2(x) \}$$

$$\pi_0(x) = \begin{cases} \text{act perform active} & x > d \\ a_3 & x = 2 \\ a_3 & x = 3 \\ a_3 & x > 4 \end{cases} \quad \begin{array}{l} (\text{forage patch 3}) \\ (\text{forage patch 3}) \\ (\text{forage patch 3}) \end{array}$$

$$M_1(x) = \begin{cases} \text{can perform action} & x=d \\ a_3 & x=2 \text{ (Forage patch 3)} \\ a_3 & x=3 \text{ (Forage patch 3)} \\ a_1 & x=4 \text{ (Forage patch 1)} \end{cases}$$

$$M_2(x) = \begin{cases} \text{can perform action} & x=d \\ a_3 & x=2 \text{ (forage patch 3)} \\ a_1 & x=3 \text{ (forage patch 1)} \\ a_1 & x=4 \text{ (forage patch 1)} \end{cases}$$

- ② a) State: Let the state at the k th stage represent the crop production at that stage.
Action: The action we perform u_k , i.e. the fraction of crops invested production in the k th stage invested for production in $(k+1)$ th stage
Space
(x_k)

State evolution: $x_{k+1} = x_k + w_k u_k x_k$

where $\{w_k\}$ is iid

$$\text{with } E(w_k) = \bar{w} > 0$$

Here we want to obtain maximise the total crop production. So,

let cost $\$ g(x_k, u_k, x_{k+1}) = (1-u_k)x_k$
 (fraction of crops stored)

& the final cost $g_N(x_N) = x_N$

$$\Rightarrow J^k = \max_{u_0, u_1, u_2, \dots, u_{N-1}} E \left(x_N + \sum_{k=0}^{N-1} (1-u_k)x_k \right)$$

Here the decision variables are u_k , $k \in \{0, 1, \dots, N-1\}$.

and the expectation is over the

sequence of random variables $\{w_k\}$

Thus by using this MDP, we have set
the appropriate objective as desired in the
problem as well as obtained an appropriate
descript.

b) We can use the DP algorithm to solve this problem

$$J_k^*(x_k) = \max_{u_k} E \left(J_{k+1}^*(x_{k+1}) \xrightarrow{\text{transition}} + (1-u_k)x_k \right)$$

$$\Leftrightarrow J_N^*(x_N) = g_N(x_N) \quad \text{expectation over } u_k \text{ or } x$$

This procedure can be used computationally to
arrive at a solution but numerically since
we don't know the structure of $J_{k+1}^*(x_{k+1})$,

I would like to evaluate the first 2 steps

& take a guess and then proceed by induction

to generalise the proof $\xrightarrow{\text{transition}} g(\cdot)$

$$J_{N-1}^*(x_{N-1}) = \max_{u_{N-1}} E \left((1-u_{N-1})x_{N-1} + J_N^*(x_N) \right)$$

$$\stackrel{\text{using}}{\Leftarrow} = \max_{w_{N-1}} E_{w_{N-1}} \left((1-w_{N-1})(x_{N-1}) + [w_{N-1} + w_{N-1}u_{N-1}x_{N-1}] \right)$$

The state equation

$$= \max_{w_{N-1}} \left[(1-w_{N-1})(x_{N-1}) + x_{N-1} + \bar{w} w_{N-1} x_{N-1} \right]$$

$$(\because E(w_k) = \bar{w})$$

Notice the eqn is linear in u_{N-1} so the extrema lie at the ends of the ~~bounds~~ allowed values of u_{N-1}

Sub $u_{N-1} = 0 \neq 1$

$$\Rightarrow T_{N-1}^*(x_{N-1}) = \max_{u_{N-1}} \{ 2x_{N-1}, (1+\bar{\omega})x_{N-1} \}$$

$$2x_{N-1} > (1+\bar{\omega})x_{N-1} \Rightarrow \bar{\omega} > 1 \quad (\because x \text{ is +ve})$$

$$\text{So, if } \bar{\omega} > 1, u_{N-1}^* = 1$$

$$\text{if } \bar{\omega} \leq 1, u_{N-1}^* = 0$$

& the $T_{N-1}^*(x_{N-1})$ are $\begin{cases} (1+\bar{\omega})x_{N-1} & \bar{\omega} > 1 \\ 2x_{N-1} & \bar{\omega} \leq 1 \end{cases}$
 → transition catg(.) correspondingly

$$T_{N-2}^* = E \left((1 - u_{N-2}) x_{N-2} + T_{N-1}^*(x_{N-1}) \right)$$

If $\bar{\omega} > 1$,

$$T_{N-2}^*(x_{N-2}) = E \left((1 - u_{N-2})(x_{N-2}) + (1 + \bar{\omega})x_{N-1} \right)$$

$$= E \left((1 - u_{N-2})(x_{N-2}) + (1 + \bar{\omega})(x_{N-2} + \bar{\omega}u_{N-2}x_{N-2}) \right)$$

$$= (1 - u_{N-2})(x_{N-2}) + (1 + \bar{\omega}) \left(x_{N-2} + \frac{u_{N-2}x_{N-2}}{1 + \bar{\omega}} \right)$$

$$= \max \left((2 + \bar{\omega})(x_{N-2}), (x_{N-2})(1 + \bar{\omega})^2 \right)$$

Since $\bar{\omega} > 1$,

$$J^+_{N-2} = (1 + \bar{\omega})^2 x_{N-2} \quad \text{if } u_{N-2}^+ = 1$$

$$J^+_{N-k} = 1 + J^+_{N-k-1} \quad \text{if } u_{N-k}^+ = 1 + \bar{\omega}$$

(Given : If $\bar{\omega} > 1$, we have shown that for $k=1$ & $k=2$.

Now we can say showing it for $k+1$
if k is true.

$$J^+_{N-k-1} = \max_{u_{N-k-1}} \left((1 - u_{N-k-1}) x_{N-k-1} + J^+_{N-k} (x_{N-k}) \right)$$

$$= \max_{u_{N-k-1}} \left((1 - u_{N-k-1}) x_{N-k-1} + (1 + \bar{\omega})^{(x_{N-k-1} + x_{N-k})} \bar{\omega} u_{N-k-1} \right)$$

$$= \max \left\{ x_{N-k-1} \left(1 + (1 + \bar{\omega})^k \right) \middle| x_{N-k-1} \left(1 + \bar{\omega} \right)^{k+1} \right\}$$

$$= \max \left\{ x_{N-k-1} \frac{\left(1 + (1 + \bar{\omega})^k \right)}{(x_{N-k-1})(1 + \bar{\omega})^{k+1}} \right\}$$

$$k+1 \cdot (1 + \bar{\omega})^k > 1 \Rightarrow 2(1 + \bar{\omega})^k > 1 + (1 + \bar{\omega})^k$$

$$\Rightarrow (1 + \bar{\omega})^{k+1} > 1 + (1 + \bar{\omega})^k \quad \rightarrow \because 1 + \bar{\omega} > 1 + 1 = 2$$

$$= x_{N-k-1} (1 + \bar{\omega})^{k+1}$$

$$\therefore u_{N-k-1}^+ = 1$$

thus the proposition holds for all k . (by induction)

$$\text{if } \mu_{N-k}^+ = \frac{\text{if } \bar{\omega} > 1}{(\varepsilon \tau_{N-k}^+(x_{N-k}) = (1+\bar{\omega})^k x_{N-k})}$$

If $\bar{\omega} < 1$

$$\begin{aligned} \tau_{N-2}^+(x_{N-2}) &= \max_{\mu_{N-2}} E((1-\mu_{N-2})(x_{N-2}) \\ &\quad + \tau_N^+(x_{N-1})) \\ &= \max_{\mu_{N-2}} E((1-\mu_{N-2})(x_{N-2}) + 2x_{N-1}) \\ &= \max_{\mu_{N-2}} E((1-\mu_{N-2})(x_{N-2}) + 2 \left(\frac{x_{N-2} + \bar{\omega} x_{N-1}}{\mu_{N-2} w_{N-2}} \right)) \\ &= \max \{ 3x_{N-2}, (2+\bar{\omega})x_{N-2} \} \\ \therefore \bar{\omega} < 1, \quad 3 &\cancel{>} 2+\bar{\omega} \end{aligned}$$

$$\Rightarrow \tau_{N-2}^+(x_{N-2}) = 2x_{N-2}.$$

$$\text{Claim: } \tau_{N-k}^+(x_{N-k}) = k x_{N-k} \quad \text{if } \bar{\omega} < 1$$

$$\text{if } \mu_{N-k}^+ > 0$$

$$\begin{aligned} \tau_{N-k-1}^+ &= \max_{\mu_{N-k-1}} E((1-\mu_{N-k-1})x_{N-k-1} \\ &\quad + k(x_{N-k-1} + \frac{\bar{\omega} x_{N-k}}{\mu_{N-k-1}})) \\ &= \max_{\mu_{N-k-1}} \left[(1-\mu_{N-k-1})x_{N-k-1} \right. \\ &\quad \left. + k(x_{N-k-1} + \mu_{N-k-1}(\bar{\omega} x_{N-k})) \right] \end{aligned}$$

(3)

$$= \max \left\{ (k+1) \bar{\omega} x_{N-k-1}, x_{N-k-1} (\bar{\omega}) (1+\bar{\omega}) \right\}$$

$$(k+1) x_{N-k-1} < k x_{N-k-1} (1+\bar{\omega})$$

$$\Rightarrow (1+\bar{\omega}) < \frac{k+1}{k} \Rightarrow \bar{\omega} < \frac{1}{k}.$$

So if $\bar{\omega} \leq \frac{1}{N}$, we can complete
the induction & the claim modified claim!

$$J^+_{N-k}(x_{N-k}) = k x_{N-k}, \text{ if } u_{N-k} = 0 \quad \forall k$$

$$\text{if } \bar{\omega} \leq \frac{1}{N} \text{ holds}$$

However if $\bar{\omega} < 1$ but $\bar{\omega} > \frac{1}{k}$,

new claim: choose $u_{N-j}^k = 0 \quad \forall j \in \{1, \dots, k\}$
 $u_{N-j} = 1 \quad \forall j > k$

where k is such that

$$\frac{1}{k+1} \leq \bar{\omega} \leq \frac{1}{k}$$

From the previous proof we can infer that

$$\text{as long as } \bar{\omega} < \frac{1}{j}, \quad u_{N-j}^k = 0.$$

We still need to show the final part: $u_{N-j}^k \geq \frac{1}{\bar{\omega}} > \frac{1}{j}$

So consider a part after $N-k+1$

say $N-k-j$, $j > 0$

we have shown that $u_{N-k-1}^* = 1$
(base case)

Now for if we can try induction

Assume $T_{N-k-j}^k = (k)(1+\bar{\omega})^j x_{N-k-j}$

($\because \mu \rightarrow 0$ being followed
by $u=1$ j times)

$$T_{N-k-j-1}^k = \max_{u \in \mathbb{R}} \left((1-u_{N-k-j-1}) x_{N-k-j-1} + T_{N-k-j}^k \right)$$

$$= \max_{u \in \mathbb{R}} \left((1-u_{N-k-j-1}) x_{N-k-j-1} + k(1+\bar{\omega})^j (x_{N-k-j-1} + u_{N-k-j-1}) \right)$$

$$= \max_{u \in \mathbb{R}} \left[(1-u_{N-k-j-1}) + k(1+\bar{\omega})^j (x_{N-k-j-1} + \bar{\omega} u_{N-k-j-1}) \right]$$

$$= \max \left\{ x_{N-k-j} \left(1 + k(1+\bar{\omega})^j \right), x_{N-k-j} \left(1 + k(1+\bar{\omega})^{j+1} \right) \right\}$$

$$\begin{aligned} &\downarrow \\ &u=0 \\ &(1-\bar{\omega})^{j+1} \\ &\text{vs } \downarrow \\ &u=1 \\ &(1+k(1+\bar{\omega})^j) \end{aligned}$$

$$\Rightarrow \max \left(\left[(1+\bar{\omega})^j k \right] \{ 1 + \bar{\omega} \}, \left[(1+\bar{\omega})^{j+k} \right] \left(\frac{1}{k(1+\bar{\omega})} + 1 \right) \right)$$

removing common factors

$$\cdot \max \left((1+\bar{\omega}), \frac{1}{k(1+\bar{\omega})} + 1 \right)$$

$$\bar{\omega} > \frac{1}{k} \Rightarrow \bar{\omega} > \frac{1}{k+1+\bar{\omega}}$$

$$\therefore T_{N-k-j-1}^+ = (1+\bar{\omega})^{j+k} \quad \text{by same} \\ \text{factors} > 1$$

$\mathcal{U} = 1$ better.

By principle of mathematical induction

our claim is true

Optimal policy

$$\text{i) If } \bar{\omega} > 1, u_0^* = u_1^* = \dots = u_{N-1}^* = 1$$

$$\text{ii) If } 0 < \bar{\omega} < \frac{1}{N}, u_0^* = u_1^* = \dots = u_{N-1}^* = 0$$

$$\text{iii) If } \frac{1}{N} < \bar{\omega} < 1, u_0^* = u_1^* = \dots = u_{N-k-1}^* = 1$$

$$u_{N-k}^* = u_{N-k+1}^* = \dots = u_{N-1}^* = 0$$

where k is such that $\frac{1}{k+1} < \bar{\omega} \leq \frac{1}{k}$

③ a) For an optimal policy π^* having cost J^*

$$T_{\pi^*} J^* = TJ^* \Rightarrow J^* = TJ^* \quad (\text{prop 2})$$

and J^* is the unique solution to the above equation ①

additionally, J^* is the unique solution to

$$T_{\pi^*} J^* = J \quad \text{--- ②}$$

Upon introducing a new action a' in a

state s' , we get a new MDP. Let the optimal policy of such an MDP be π^{new} and

cost be J^{new} .

From ②, since we don't have a change in transition probabilities & π^* (i) is also unchanged

$$T_{\pi^*} J^* = J^* \text{ still holds} \quad \text{--- ③}$$

We still need to check if

$J^* - TJ^*$ to know whether J^* is indeed J^{new} .

[Assuming minimisation problem for part a)
- although it really doesn't matter here]

$\therefore \pi^*$ was optimal for the older MDP

$$(T\pi^*)(i) = \pi^*(i) \text{ in older MDP}$$

$$\Rightarrow \pi^*(i) = \min_a \sum_{j \in S} (P_{i,j}(a) (g(i,a,j) + \pi^*(j)))$$

In the new MDP,

$$T\pi^*(i) = \min_a \sum (P_{i,j}(a) (g(i,a,j) + \pi^*(j)))$$

holds for all $i \neq s'$ — (4)

(because nothing has changed
— no new actions have been added)

claim(1): If $T^*(s') = \min_a \sum_{j \in S} (P_{s',j}(\pi^*(s')) (g(s',\pi^*(s'),j) + \pi^*(j)))$

$$= \min_a \sum (P_{s',j}(a) (g(s',a,j) + \pi^*(j)))$$

then, $\pi^* = \pi^*$ optimal & $T^*(s') = \pi^*$ optimal

Proof: Procedure to test :

Assume that compare $T^*(s')$ with

$$\sum_{j \in S} P_{s',j}(a') (g(s',a',j) + T^*(j))$$

$$\text{if } T^*(s') \leq \sum_{j \in S} P_{s',j}(a') (g(s',a',j) + T^*(j))$$

(2)

then $\pi^*(s)$ is still the solution to the eqn because substituting $\pi^*(s)$ gives us a lower cost than a^*

$\Rightarrow \pi^*(s)$ is still optimal

& this is the claim is satisfied.

$$\Rightarrow \pi^* = \pi^{\text{new}} \quad \pi^* = \pi^{\text{new}}$$

Proof for the claim:

$$J^*(s) = \min_a \sum p_{s'j}(a) (g(s', a, j) + J^*(s'))$$

then combining with eqn ④, we have
(for new MDP),

$J^* = J^*$ $\Rightarrow J^*$ is still the optimal cost
as we already showed that
expected cost of policy π^* is J^*

$$\Rightarrow J^*_{\text{new}} = J^* \quad \& \quad \pi^{\text{new}} = \pi^*$$

Summary: For the s_1 state equation, check if action a_1 produced a minimum ^{in RHS} if it doesn't the old policy remains the optimal policy

claim 2: If a' gives a lower value in $\pi^* + \gamma^*$ (note that γ^* is the optimum in old MDP, it is not so in the new MDP)

for s' state; then for that state s', π^* , action a' is optimal.

Proof:

Suppose a' is not the optimal action, then it is as good as dealing with an MDP without that action.

\Rightarrow this is the old MDP

$\Rightarrow \pi^* \in \gamma^*$ are optimal in the new MDP then

$$\Rightarrow \gamma^*(s') = \min_a \sum_{\pi_j} p_{\pi_j} (g(s', a, j) + \gamma^*(j))$$

has taken $a \in \pi^*$ instead of a'
 this is a contradiction because $\gamma^*(s) \neq \sum p_{\pi_j} (\cdot)$
to find new policy \Rightarrow the claim is true.

From claim 2 we can infer that a' is the optimal action at s' if $\pi^* \neq \pi^*$ new

a) Let $\pi_0^*(i)$ s.t $\pi^*(i) = \begin{cases} \pi^*(i) & i \neq s^* \\ a' & i = s^* \end{cases}$

Using this policy we can perform policy iteration to find the new optimal policy

[note that $\pi_0(i)$ is proper since

$$J^* \geq T_{\pi_0} J^* - (\text{if } a' \text{ is optimal})$$

as shown in that eqn in claim ①]

Since we already found one optimal action
policy iteration might converge faster

b) Assuming it is a maximization problem
we are given that the cumulative reward

$J^*(s^*)$ decreases by Δ

$$J_{\pi^*}(s^*) = E \left(g(s^*, \pi^*(s^*), j) + J^*(j) \right)$$

Since no other information is given, one π^*
 would like to interpret it as $T_{\pi^*}(s^*)$
 has changed because of a Δ decrease in
 average expected transition^{cost} from s^*

$$\text{that } \pi^* \in (g_{\text{new}}(s^*, \pi^*(s^*), j)) \\ = E(g_{\text{new}}(s^*, \pi^*(s^*), j))$$

Because of this change in transition cost from s^* ,
 there are 2 implications

- i) T_{π^*} might not be solution to $\pi^* = T_{\pi^*}$
 $T = T_{\pi^*} T$
 i.e. it may not be the expected cost of policy π^*
- ii) π^* may not be optimal policy (because
 return from s^* has reduced, intuitively
 there is a possibility that a policy that
 return has lesser number of transitions to s^*
 may be the new optimal policy, π^{new})

Checking if π^* is optimal

Step ①: Perform policy evaluation

$$J(\pi) = \sum_j p(j|\pi(i)) g(i, \pi_k^*(i), j) + J(j))$$

where $g(i, \pi_k^*(i), j) = \begin{cases} g_{\text{original}}(i, \pi_k^*(i), j) & \text{if } i \in S^* \\ 0 & \text{otherwise} \end{cases}$

$$J(i) = E(g(i, \pi_k^*(i), j) + J(j)) \quad \text{if } i \in S^*$$

These ∞ for $i \in S^*$

$$J(S^*) = E(g_{\text{original}}(S^*, \pi_k^*(S^*), j) + J(S^*)) - \Delta$$

where g_{original} is the transition prob

of original MDP. Let's solve these equations.

Step ②: Perform policy improvement:

$$\pi_{\text{new}} J^* = T J^*$$
$$\pi_{\text{new}} J^* = \max_a E(g(i, a, j) + J(j))$$

If J_{new}

Get π_{new} & evaluate J_{new} .

If $J_{\text{new}} = J^*$ then π^* is the optimal policy

If $J_{\text{new}} < J^*$, π^* is not optimal.

In order to converge to the new optimal policy we can do policy iteration with J_{new} as the initial proper policy.

The expected cost will also be obtained via policy iteration in policy evaluation by step

(4)

For convenience, let's consider Kevin as player A, Karen as player B, Kannan as player C.

To model the problem as an MDP, we will assume that 'A' is the decision maker and we need to know the decisions (in fact, optimal) of 'B' & 'C'.

2-Game scenarios .

Consider the game has reduced to A & B and it is A's turn to fire. Let us analyse this from the point of view of B:

- o B can't miss intentionally - has to shoot at A.
- o Rule permits A to shoot into the air. But that will increase chance of B winning - because B will take more turns to fire so hence the probability of reaching the target (which will be a geometric distribution) increases

So, to obtain a lower bound of probability of B winning we can as well assume that A always shoots & doesn't miss intentionally.

$$P_{\text{B wins in } k \text{ turns}} = P(\text{A misses } k \text{ turns despite shooting}) \\ \times P(\text{B misses } k-1 \text{ turns correctly before firing its last turn})$$

(1)

$$= (1-\alpha)^k (1-\beta)^{k-1} \beta.$$

$$\therefore P(B \text{ wins}) = \sum_{k=1}^{\infty} (1-\alpha)^k (1-\beta)^{k-1} \beta.$$

~~if A vs B~~

$$= \frac{(1-\alpha)^k \beta}{1 - (1-\alpha)(1-\beta)} \times \beta$$

Proceeding similarly

$$P(B \text{ wins in } B \text{ vs } C) = \frac{(1-\alpha)r}{1 - (1-r)(1-\beta)} \beta$$

$$(1-r) < (1-\alpha) \quad (\because \alpha > r)$$

$$\left(\Rightarrow 1 - (1-\alpha)(1-\beta) < 1 - (1-r)(1-\beta) \right)$$

So we conclude

$$P(B \text{ wins } A \text{ vs } B) > P(B \text{ wins } B \text{ vs } C)$$

Similarly for C arguing in a similar fashion,

$$P(C \text{ wins } C \text{ vs } A) > P(C \text{ wins } C \text{ vs } B)$$

From this analysis, we understand that

B and C would seek to eliminate each other first

because they prefer duelling A in a $\frac{2}{3}$ game scenario.
to minimize winning chance

So, when all of ABC are alive,

* B shoots C in his turn

* C shoots B in his turn

Bottom

Now that we have decided the possible actions of B & C we can build an MDP.

- Note that a single transition in the MDP to incorporate all changes in the position of the game between 2 turns of A (eg. if B & C are alive, the transition implies A, B, and C have all taken their turns)

State Space : {
 ① A ~~dead~~, A and B are alive, A and C are alive,
 ② A and B are alive, A and C are alive,
 ③ A, B and C are alive, A only alive (A has won!) }
 ④ A, B and C are alive, A only alive (A has won!) }

Action Space :

state ① : no action - A already out of the game (terminal)

state ② : shoot at B, miss intentionally.

state ③ : shoot at C, miss intentionally

state ④ : shoot at B, shoot at C, miss intentionally

state ⑤ : A has won! - no action (terminal)

Cost/Rewards:

We want to minimise probability of winning,

$$g(i, a|j) = 0 + \{ 1, 0 \}$$

~~we will normalise~~ $g(i, a|j) = 0 \text{ if } x \neq x_5 \\ 1 \text{ if } x = x_5$

Note that the $G_N(x_N) = \{ 0 \text{ if } x \neq x_5 \\ 1 \text{ if } x = x_5 \}$
 A last 4 turns are terminal states

$$\begin{aligned}
 & \text{So objective max } \underset{a}{\mathbb{E}} \left(\alpha + J_N^+(x_N) + \sum_{k=1}^N g_k(x_k, u_k, u_{k+1}) \right) \\
 & = \max_a \mathbb{E}(J_N^+(x_N)) \\
 & = \max_a \mathbb{E}(g_N(x_N))
 \end{aligned}$$

which is eqvt to minimizing probability of miss.

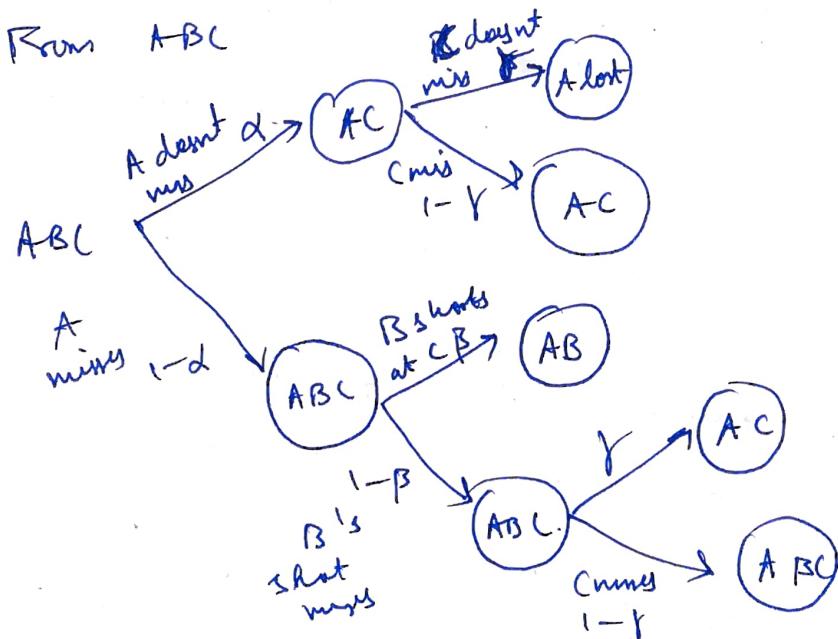
\rightarrow C: probability = expectation of indicator variable for the event)

Transition probabilities

Let a_1 : A shoots B a_2 : A shoots C as A shoots air.

a) Action a_1 :

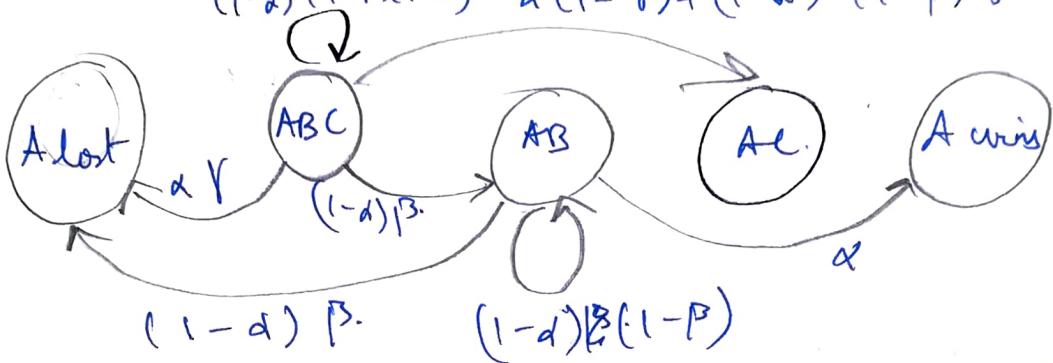
From AB: A wins if A shoots properly
 \downarrow
 & loses if A misses & B shoots properly
 Back to same state iff A & B both miss $(1-\alpha)(1-\beta)$



Using these transition probabilities we get the

markov chain as:

$$(1-\alpha)(1-\beta)(1-\gamma) \propto (1-\delta) + (1-\alpha)(1-\beta)\gamma$$



No transitions from AC because to cut shoot B in that state -

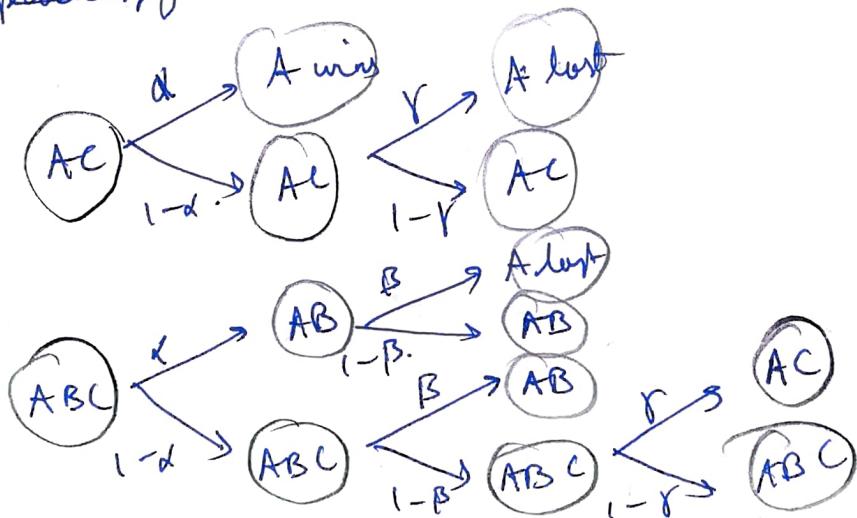
$$P_{21} = \alpha\gamma \quad P_{22} = (1-\alpha)(1-\beta)(1-\delta) \quad P_{24} = \\ \alpha(1-\delta) + (1-\alpha)(1-\beta)\gamma$$

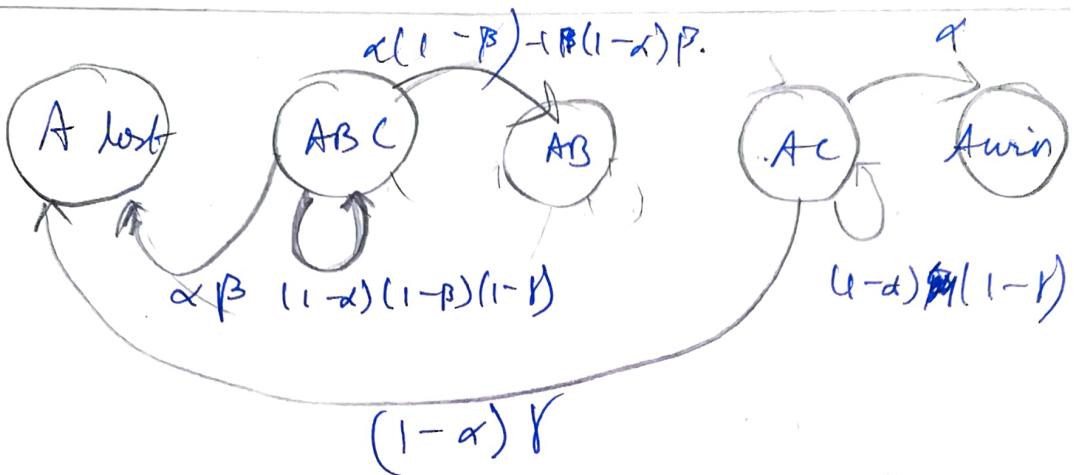
$$P_{23} = (1-\alpha)\beta$$

$$P_{31} = (1-\alpha)\beta \quad P_{33} = (1-\alpha)(1-\beta)$$

$$P_{35} = \alpha \quad ; \text{ other } P_{ij} = 0$$

ii) Action as : A shoots C
probability proceeding similarly



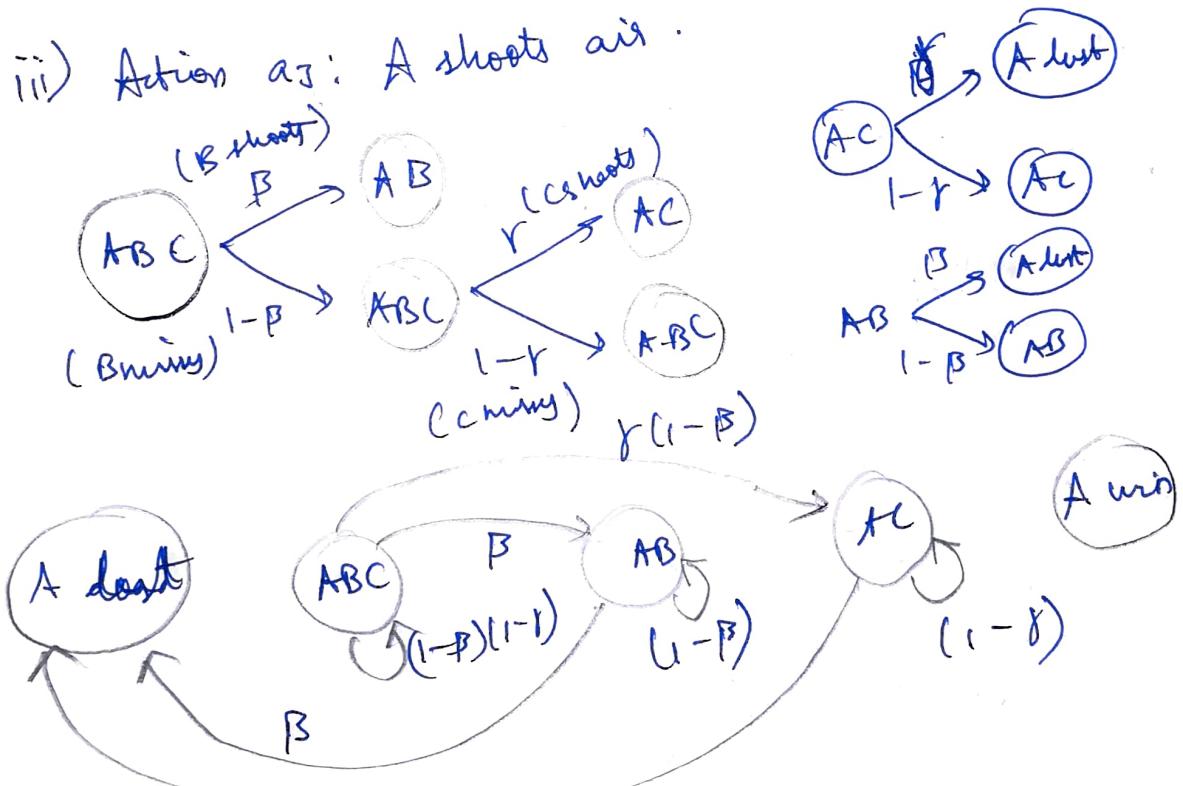


$$P_{21} = \alpha\beta \quad P_{22} = (1-\alpha)(1-\beta)(1-\gamma)$$

$$P_{23} = \alpha(1-\beta) + (1-\alpha)\beta; P_{21} = (1-\alpha)\gamma$$

$$P_{24} = (1-\alpha)\beta(1-\beta); P_{45} = \gamma \text{; other } P_{ij} = 0$$

iii) Action a_3 : A shoots air.



$$P_{22} = (1-\beta)(1-\beta) = r^2 \quad P_{33} = (1-\beta) \quad P_{44} = 1-r$$

$$P_{23} = \beta \quad P_{24} = (1-\beta)r \quad P_{31} = \beta \quad P_{41} = r$$

other $P_{ij} = 0$

b) Bellman optimality equations

$$TJ_N^+ = J^+ ; \quad J^+(A_{\text{win}}) = 1$$

$$\Rightarrow J^+(i) = \max_{a \in A(i)} E(g(Y_i, a, j) + J^+(j))$$

(- : maximizing rewards)

$$= \max_{a \in A(i)} E(J^+(j))$$

$J(A|B) = \max_a E$

At $A|B$ state $\#B$, we have 2 actions allowed
- shoot B or shoot at

$$\Rightarrow J(A|B) = \max \left\{ \begin{array}{l} (1-\beta) J(A|B), \\ \text{A shoot } B \\ \text{or} \\ \text{A shoot } B \end{array} \right. + \alpha \left. \begin{array}{l} (1-\beta) J(A|B) \\ + \alpha \end{array} \right\}$$

$$= \max \left\{ \begin{array}{l} (1-\beta) J(A|B), \\ (1-\alpha)(1-\beta) J(A|B) \\ + \alpha \end{array} \right\}$$

(1)

Similarly

$$J(A|C) = \max \left\{ \begin{array}{l} (1-\gamma) J(A|C), \\ (1-\alpha)(1-\gamma) J(A|C) \\ + \alpha \end{array} \right\}$$

(2)

(2 actions : shoot air or shoot C)

At node ABC, all 3 actions are allowed, so we take max over the 3 actions

$$\begin{aligned}
 J(ABC) = \max \{ & [(1-\alpha)(1-\beta)(1-\gamma)J(AB)] \\
 & + (\alpha(1-\gamma) + (1-\alpha)(1-\beta))J(Ac) \\
 & + (1-\alpha)\beta J(BC), \\
 & [(1-\alpha)(1-\beta)(1-\gamma)J(ABC) \\
 & + (\alpha(1-\beta) + (1-\alpha)\beta)J(AB) \\
 & + (1-\alpha)\gamma(1-\beta)J(AC)], \\
 & [(1-\beta)(1-\gamma)J(ABC) \\
 & + \beta J(AB) + \gamma(1-\beta)J(AC)] \} \\
 & \xrightarrow{\text{---}} ③
 \end{aligned}$$

Given $\alpha = 0.3$, $\beta = 0.5$, $\gamma = 0.6$. Substitute,

$$\text{eqn } ① : J_{AB} = \max (0.5J_{AB}, 0.5 \times 0.7 \\
 = 0.3J_{AB} + 0.3)$$

part i) giving $\beta + \gamma = 0$ - not correct - they must probability

$$\text{part ii) using } J_{AB}^+ = \frac{0.3}{0.65} = \boxed{0.4615} - \frac{\text{if my}}{\text{I win } J_{AB} > 0} \quad \text{if my} \\
 \text{I win } J_{AB} > 0 \\
 - \frac{\text{shoot B prefud}}{(0.4615 \times 0.5)}$$

$$\text{eqn } ② : J_{Ac}^+ = \max (0.4J_{Ac}, 0.28J(AC) + 0.3) \\
 < 0.35(0.4615) + 0.3$$

using a similar argmt,

$$J_{Ac}^+ = \frac{0.3}{0.72} = \boxed{0.4167} - \frac{\text{shoot C}}{\text{prefud}}$$

eqn ③

(i) term :

$$\begin{aligned} J_{ABC}^+ &= (0.7)(0.5)(0.4) J_{ABC} \\ &\quad + [(0.3)(0.4) + (0.7)(0.5)(0.6)] J_{AC} \\ &\quad + (0.7)(0.5) J_{AB} \end{aligned}$$

$$\Rightarrow J_{ABC}^+ = \frac{0.299}{0.86} = 0.3477$$

ii) term :

$$\begin{aligned} J_{ABC}^+ &= 0.14 J_{ABC} + [(0.3)(0.5) + (0.7)(0.4)] J_{AC} \\ &\quad + (0.5)(0.7)(0.6) J_{AB} \end{aligned}$$

$$\Rightarrow J_{ABC}^+ = \frac{0.3183}{0.86} \approx 0.37$$

iii) term :

$$\begin{aligned} J_{ABC}^+ &= 0.2 J_{AB} + 0.5 J_{AC} + 0.3 J_{BC} \\ \Rightarrow J_{ABC}^+ &= \frac{0.351}{0.8} = 0.445 \end{aligned}$$

Using $J_{AC}^+ = 0.445$,

(iii) term $>$ (i) \therefore eqn (ii)

\Rightarrow eqn is solved.

— shoot in air
preferred.

S_{start}	$f^*(x)$	$\text{optimal } f^*(x)$ (optimal action)
A not	0	— (A can do anything)
AB	0.4615	Shoot B
AC	0.4667	Shoot C
ABC	0.445	Shoot in air

⑤ Before I answer the question, I would like to show $\|a + b\|_g < \|a\|_g + \|b\|_g$

$$\|a + b\|_g = \max_i \frac{|a(i) + b(i)|}{\varepsilon(i)}$$

$$[\text{triangle inequality}] \leq \max_i \left(\frac{|a(i)| + |b(i)|}{\varepsilon(i)} \right)$$

$$\because \max((a+b)(i)) \leq \max_i \frac{|a(i)|}{\varepsilon(i)} + \max_j \frac{|b(j)|}{\varepsilon(j)}$$

$$\therefore = \|a\|_g + \|b\|_g$$

$$\Rightarrow \|a + b\|_g < \|a\|_g + \|b\|_g$$

_____ ①

a) Termination criterion: $\|T_{m+1} - T_m\|_g < \varepsilon \left(\frac{1-\beta}{2\beta} \right)$

_____ ②

T is a contraction operator

$$\Rightarrow \|TJ - TJ'\|_g \leq \beta \|J - J'\|_g$$

_____ ③

Apply operator T in eqn ②

$$\Rightarrow \|TJ_{m+1} - TJ_m\|_g < \beta \varepsilon \left(\frac{1-\beta}{2\beta} \right)$$

(- JT is a
contractive
operator)
- eqn ③

$$\Rightarrow \| T_{m+2} - T_{m+1} \|_g < \delta \left[\varepsilon \left(\frac{1-\delta}{2\delta} \right) \right]$$

Reversely applying T ,

$$\| T_{m+3} - T_{m+2} \|_g < \delta^2 \left[\varepsilon \left(\frac{1-\delta}{2\delta} \right) \right]$$

⋮

$$\| T_{m+n} - T_{m+n-1} \|_g < \delta^{n-1} \left[\varepsilon \left(\frac{1-\delta}{2\delta} \right) \right]$$

Add all these eqns.

$$\Rightarrow \sum_{k=2}^n \| T_{m+k} - T_{m+k-1} \|_g < \sum_{k=1}^{n-1} \delta^k \left[\varepsilon \left(\frac{1-\delta}{2\delta} \right) \right] \quad (4)$$

Now we use $\| a + b \|_g < \| a \|_g + \| b \|_g$

and push the sum inside the modulus

$$\begin{aligned} \text{i.e. } & \| T_{m+n} - T_{m+n-1} + T_{m+n-1} - T_{m+n-2} \\ & + \dots + T_{m+3} - T_{m+2} + T_{m+2} - T_{m+1} \|_g \\ & < \sum_{k=2}^n \| T_{m+k} - T_{m+k-1} \|_g \end{aligned}$$

$$\textcircled{4} \Rightarrow \|J_{m+n} - J_{m+1}\|_g < \sum_{k=1}^{n-1} g^k \left(\varepsilon \left(\frac{1-g}{2g} \right) \right)$$

$\not\rightarrow$ take limit
 $n \rightarrow \infty$

$$\lim_{n \rightarrow \infty} \|J_{m+n} - J_{m+1}\|_g < \lim_{n \rightarrow \infty} \sum_{k=1}^{n-1} g^k \left(\varepsilon \left(\frac{1-g}{2g} \right) \right)$$

$$\varepsilon_g = \frac{g}{1-g} \varepsilon \left(\frac{1-g}{2g} \right)$$

$$\not\rightarrow \lim_{n \rightarrow \infty} \|J_{m+n} - J_{m+1}\|_g^2 \leq \frac{\varepsilon}{2}$$

Also, we know that value iterations converge

$$\left(\lim_{k \rightarrow \infty} T^k J_i = J^* \right)$$

$$\Rightarrow \|J_m^* - J_{m+1}\|_g < \frac{\varepsilon}{2}$$

$$\Rightarrow \|J_{m+1} - J_m^*\|_g < \frac{\varepsilon}{2}$$
(5)

Hence proved.

b) From part (a) we have

$$\|J_{m+1} - J_m^*\|_g < \frac{\varepsilon}{2}$$

Applying T operator & using T is a contraction

$$\|T_{Jm+1} - J^*\| < \frac{\delta \varepsilon}{2} \quad (\because T^+ = J^*)$$

$$\text{But } T_{\pi^\varepsilon} J_{m+1} = T_{Jm+1} \quad (\because T^+ = J^*)$$

$$\Rightarrow \|T_{\pi^\varepsilon} J_{m+1} - J^+\| < \frac{\delta \varepsilon}{2} \quad \textcircled{6}$$

In showing each policy iteration indeed gives us a policy improvement,

(which is what we had done here)

we showed (in class)

$$T_{\pi^\varepsilon} \leq T_{m+1} \Rightarrow T_{\pi^\varepsilon} T_{\pi^\varepsilon} \leq T_{\pi^\varepsilon} J_{m+1}$$

\downarrow
 $(\because T_{\pi^\varepsilon} \text{ is a monotone operator})$

$$\Rightarrow T_{\pi^\varepsilon} T_{\pi^\varepsilon} - J^+ \leq T_{\pi^\varepsilon} J_{m+1} - J^+$$

$$\text{But } T_{\pi^\varepsilon} J_{m+1} = T_{Jm+1} \Rightarrow T_{\pi^\varepsilon} - J^+ \leq T_{Jm+1} - J^+$$

$$\Rightarrow \|T_{\pi^\varepsilon} - J^+\| \leq \|T_{Jm+1} - J^+\| \quad \textcircled{7}$$

$$\text{using } \textcircled{7} \text{ in } \textcircled{6} \quad \|T_{\pi^\varepsilon} - J^+\| \leq \|T_{\pi^\varepsilon} J_{m+1} - J^+\| \leq \frac{\delta \varepsilon}{2}$$

$$\Rightarrow T_{\pi^\varepsilon} \|T_{\pi^\varepsilon} - J^+\| \leq \|T_{\pi^\varepsilon} J_{m+1} - J^+\| \leq \frac{\delta \varepsilon}{2}$$

$$\Rightarrow \|T_{\pi^\varepsilon} - J^+\| \leq \frac{\delta \varepsilon}{2}; \text{ But } 0 < \delta < 1 \Rightarrow \frac{\delta \varepsilon}{2} < \varepsilon$$

$$\Rightarrow \|T_{\pi^\varepsilon} - J^+\| < \varepsilon$$

Hence proved