

Declaration :

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

Acknowledgement :

- ① Beftsekas DPOC Vol 1 Chap 1 (MOP)
- ② Class notes

① We know that finite-horizon MDPs can be converted to an equivalent SSP

Consider an MDP with state space X , N stages & action space A

Let the states in the SSP be

$$(i, k) \xrightarrow{a} i \text{th state in MDP} \\ i \in X$$

$$k \rightarrow k\text{th stage in MDP}$$

$$k \in \{1, \dots, N\}$$

Terminal states : $(i, N) \forall i \in X$

[Because finite-horizon MDP gets over after N stages]

Transitions : $(i, k) \xrightarrow{a} (j, k+1)$ with probability $p_{ij}(a)$

$$a \in A_k(i)$$

where $A_k(i)$ is the set of available actions

in the MDP at stage k in state i

Costs : $g_{SSP}(i, N), a, (j, N)) = 0$ (\because terminal)

$$g_{SSP}(i, k), a, (j, k+1)) = g_{MDP}(i, a, j)$$

Now, we can perform our usual policy iteration on the SSP

One more thing to note is All Policies are Proper in this SSP.

Proof: Note that all transitions lead from $(i, k) \rightarrow (j, k+1)$ i.e. the second index denoting the stage keeps increasing till it reaches N .

→ irrespective of policy, the terminal state will be reached in N steps.

[this is so, because the finite horizon MDP finishes in N stages]

Policy Iteration

- ① Start with a proper policy π_0 (all policies are proper so start with any policy).
- ② Policy evaluation:
Solve $J = T_{\pi_k}^* J_k$; J_{k+1} will be unique because of SS + proper

$$\Leftrightarrow J_{\pi_k}(i, r) = \sum_{\pi_j} P((i, r), \pi_k(i, r), (j, r+1)) \times \\ (g_{\text{SSP}}((i, r), \pi_k(i, r), (j, r+1)) + J_{\pi_k}(j, r+1))$$

$\Leftrightarrow P((i, r), \pi_k(i, r), (j, r+1))$ is same as

$P(i, \pi_k(i, r), j)$
MDP in the r^{th} stage

$$g_{\text{SSP}}((i, r), \pi_k(i, r), (j, r+1)) = g_{\text{MDP}}(i, \pi_k(i, r), j)$$

in the r^{th} stage

So we know P, π_k, g .

Just solve for J_{π_k}

(iii) Policy Improvement

$$T_{k+1} J_{\pi_k} = J_{\pi_k}$$

$$\text{i.e. } \pi_{k+1}(i, r) = \arg \min_{a \in A_r(i)} \sum P((i, r), \pi_k(i, r), (j, r+1)) \times \\ (g((i, r), \pi_k(i, r), j, r+1) + J_{\pi_k}(j, r+1))$$

$A_r(i)$ denotes set of
all allowed actions in MDP
at stage r

$$+ \frac{J_{\pi_k}(j, r+1)}{J_{\pi_k}(j, r+1)}$$

P, g, π_k are same defined in same way as policy evaluation step.

J_{π_k} is obtained from policy evaluation step

$$(i) \quad \text{If } J_{\pi_{k+1}}((i, r)) < J_{\pi_k}((i, r))$$

for at least one
state (i, r) ,

~~continue~~ ~~if~~ (go to

go to step (ii) & repeat.

Else stop.

Time complexity analysis

I. DP

$$J_k^+(i) = \min_a \left(g \sum_j P_{ij}(a) (g^{(i, a)} + J_k^+(j)) \right)$$

$j \in X \rightarrow$ totally n states.

To find the optimal action (a value) of each state at stage k , $n \times m$ operations required

[n multiplications, perform for m actions,
choose minimum resulting value]

d) For each stage k , $n \times (nm) = n^2 m$

(\because n states at each stage)

For N stages, $Nn^2 m$

Computational requirements of DP: $\boxed{O(Nn^2 m)}$

II. Policy Iterations

a. Policy evaluation - required per update,
for state $i(r)$, n multiplications
needed.

\Rightarrow totally : $n \times N$ states

$\therefore Nn^2$ operations / policy evaluate update.

b. Policy update improvement

For each state $i(r)$, perform Nn
multiplications for m actions. = Nm .

\Rightarrow for all states : $(Nn)(Nm)$
 $= (Nm)^2$ / iteration

(multiplications only for n states because

$(i, r\text{row}) \rightarrow (j, r+1)$ for other x stages the probability will be 0)

a. Policy iteration evaluation

Solve a system of linear equations
with N^n variables

$\Rightarrow (N^n)^3$ operations / iterations

c. Tabed policies : m actions for each
state in each stage

$$= (m)^{N^n}$$

In PI worst case, we iterate over all
policies

Worst case req'd : $O((m)((N^n)^m + (N^n)^3))$

So DP is computationally less intensive
than PI for finite horizon MDPs.

(2)

$n+1$ states: $\{0\} \cup \{1, \dots, n\}$

↙
Terminal state

$$R = \sum_{t=0}^{T-1} \gamma^t r(x_t), \quad \gamma \in (0, 1)$$

Some after T time steps, the terminal state is reached, $r(x_T) = r(x_{T+1}) = \dots = \lim_{k \rightarrow \infty} r(x_k) = \theta$

So we can add all those terms to the summation in $R \rightarrow$ it won't make any difference since all those terms are 0

$$\Rightarrow R = \sum_{t=0}^{T-1} \gamma^t r(x_t) + 0$$

$$= \sum_{t=0}^{T-1} \gamma^t (r(x_t)) + \sum_{t=T}^{\infty} \gamma^t (0)$$

$$= \sum_{t=0}^{T-1} \gamma^t r(x_t) + \sum_{t=T}^{\infty} \gamma^t \underset{\substack{\lim \\ t \rightarrow \infty}}{r(0)} \quad \because \gamma < 1$$

(4)

$$= \sum_{t=0}^{T-1} \gamma^t r(x_t) + \sum_{t=T}^{\infty} \gamma^t r(x_t)$$

$$r(x_t) \geq 0 \quad \forall t \geq T$$

$$\Rightarrow R = \frac{\sum_{t=0}^{\infty} \gamma^t r(x_t)}{1 - \gamma}$$

[did this manipulation because, otherwise
the sum ends at a random T]

$$a) J(x) = E(R|x_0=x)$$

$$= E\left(\sum_{t=0}^{\infty} \gamma^t r(x_t) \mid x_0=x\right)$$

$$< E\left(r(x_0) + \sum_{t=1}^{\infty} \gamma^t r(x_t) \mid x_0=x\right)$$

$$= r(x) + E\left(\sum_{t=1}^{\infty} \gamma^t r(x_t) \mid x_0=x\right)$$

$$(\because E(r(x_t) \mid x_0=x) = r(x))$$

$$= r + E\left(\sum_{t=1}^{\infty} \gamma^{t-1} r(x_{t+1}) \mid x_0=x\right) + r(x)$$

$$= \gamma E \left(E \left(\sum_{t=1}^{\infty} \gamma^{t-1} r(x_t) \mid x_0 = u, x_1 = x_1 \right) \right)$$

↓
over x_1 ↓ over x_2, x_3, \dots

etc...)

$$= r(x) + \gamma E \left[\underset{\text{weights}}{E \left(\sum_{t=0}^{\infty} \gamma^t r(x_{t+1}) \mid x_0 = x_1, x_1 = x_1 \right)} \right]$$

(re-ordering → ①)

note that $P(x_2 \mid x_1 = x_1, x_0 = u)$ [the sum]

$$= P(x_2 \mid x_1 = x_1)$$

i.e. future transitions from x_1 don't depend on x_0
[Markov property]

So the sequence of states

$\{x_0, x_1, x_2, \dots\}$ depend only on

~~at time t FPE~~ can rewrite it as $\{x_1, x_2, \dots\}$

~~$E \left(\sum_{t=0}^{\infty} \gamma^t r(x_{t+1}) \mid x_0 = x_1, x_1 = x_1 \right)$~~

$$= E \left(\sum_{t=0}^{\infty} \gamma^t r(x_{t+1}) \mid x_1 = x_1 \right)$$

can rewrite $\{x_2, x_3, x_4, \dots\}$ as

$$\{x'_1, x'_2, x'_3, \dots\}$$

where transition to x'_1 (equally x_2) is governed by x'_0 (equally $x' \rightarrow \cdot : x_1 = x'$)

$$\therefore E \left(\sum_{t=0}^{\infty} \gamma^t r(x_{t+1}) \mid x_0 = x, x_1 = x' \right)$$

$$= E \left(\sum_{t=0}^{\infty} \gamma^t r(x_{t+1}) \mid x_1 = x' \right)$$

$$= E \left(\sum_{t=0}^{\infty} \gamma^t r(x'_t) \mid x'_0 = x' \right)$$

$$= E \left(\sum_{t=0}^{\infty} \gamma^t r(x'_t) \mid x'_0 = x' \right)$$

$$= g(x') E(R \mid x'_0 = x')$$

$$= J(x') - ②$$

$$\therefore J(x) = r(x) + \gamma E(J(x'_1) \mid x'_0 = x)$$

$$J(x) = r(x) + \gamma \sum_{x' \in X} P(x' \mid x) J(x')$$

③

⑥

$$\begin{aligned}
 \mathbb{E} N(x) &= E(R^2 | x_0 = y) \\
 &= E\left(\left(\sum_{t=0}^{\infty} \gamma^t r(x_t)\right)^2 | x_0 = y\right) \\
 &= E\left(\left(r(x) + \sum_{t=1}^{\infty} \gamma^t r(x_t)\right)^2\right) \\
 &= E\left(r^2(x) + \left(\sum_{t=1}^{\infty} \gamma^t r(x_t)\right)^2 + 2r(x) \sum_{t=1}^{\infty} \gamma^t r(x_t)\right) \\
 &= r^2(x) + 2r(x) E\left(\sum_{t=1}^{\infty} \gamma^t r(x_t) | x_0 = y\right) \\
 &\quad + E\left(\left(\sum_{t=1}^{\infty} \gamma^t r(x_t)\right)^2 | x_0 = y\right) \quad \text{--- (4)}
 \end{aligned}$$

expectation is
the second term was shown to be equal

$$+ E(J(x_1^*) | x_0 = y).$$

\downarrow
on x_1^*

And by using markov property

i.e. using the equation of

$$Q(x)P(\{x_2, x_3, \dots\} | x_0, x_1, x_1)$$

$$= P(\{x_2, x_3, \dots\} | x_1)$$

$$E \left(\left(\sum_{t=1}^{\infty} r^t r(x_t) \right)^2 \mid x_0 = x \right)$$

$$= \gamma^2 E \left(E \left(\left(\sum_{t=1}^{\infty} r^{t-1} r(x_t) \right) \mid x_0 = x \right) \mid x_1 = x' \right)$$

$$= \gamma^2 \underset{\text{over } x_1}{\downarrow} E \left(E \left(\left(\sum_{t=1}^{\infty} r^{t-1} r(x_t) \right) \mid x_0 = x \right) \mid x_1 = x' \right)$$

$$= \gamma^2 E \left((E \left(\left(\sum_{t=0}^{\infty} r^t r(x_{t+1}) \right) \mid x_1 = x' \right)) \mid x_0 = x \right)$$

$$= \gamma^2 E \left((E \left(\left(\sum_{t=0}^{\infty} r^t r(x'_t) \right) \mid x'_0 = x' \right)) \mid x_0 = x \right)$$

$$= \gamma^2 E \left(M(x') \mid x_0 = x \right)$$

$$= \gamma^2 E \left(M(x') \mid x_0 = x \right) - \textcircled{5}$$

$$\textcircled{4} \ni M(x) = r^2(x) + 2\gamma r(x) E(J(x') \mid x_0 = x) + \gamma^2 E(M(x') \mid x_0 = x)$$

$$\begin{aligned} M(x) &= r^2(x) + 2\gamma r(x) \sum_{x' \in X} P(x' \mid x) J(x') \\ &\quad + \gamma^2 \sum_{x' \in X} P(x' \mid x) M(x') \end{aligned}$$

\textcircled{6}

$$\begin{aligned}
 b) \quad \text{var } V(x) &= \text{var } (R|x_0=x) \\
 &= E(R^2|x_0=x) - (E(R|x_0=x))^2 \\
 &= M(x) - (J(x))^2 \quad \text{--- (7)} \\
 &\text{use (3) \& (6) here to subst for } M(x) \text{ \&} J(x)
 \end{aligned}$$

$$\begin{aligned}
 &= r^2(x) + 2\gamma r(x) \sum p(x'|x) J(x') \\
 &\quad + \gamma^2 \sum p(x'|x) M(x') - (r(x) + \gamma \sum p(x'|x) J(x')) \\
 &\stackrel{\text{using } (a+b)^2 = a^2 + b^2 + \text{cross terms}}{\approx} r^2 \sum p(x'|x) M(x') - \gamma^2 \left(\sum p(x'|x) J(x') \right)^2
 \end{aligned}$$

Add and subtract $\gamma^2 \sum p(x'(y)) (J(y))^2$

$$\begin{aligned}
 &= \gamma^2 \left(\sum_{x' \in X} p(x'|x) J(x')^2 - \left(\sum_{y \in X} p(y|x) J(y) \right)^2 \right) \\
 &\quad + \gamma^2 \left[\sum_{x' \in X} p(x'|x) \left(M(x') - (J(x'))^2 \right) \right]
 \end{aligned}$$

$$\begin{aligned}
 &= \gamma^2 \psi(x) + \gamma^2 \sum_{x' \in X} p(x'|x) V(x') \\
 &\quad \cdot \left(\begin{array}{l} \text{from eqn 7 } \text{var}(y) \\ = M(y) - (J(y))^2 \\ \text{from definition of } \psi(y) \end{array} \right)
 \end{aligned}$$

$$(3) \text{ (a) contraction property } \|f(a) - f(b)\|_{\infty} \leq \beta \|a - b\|_{\infty}$$

$$\text{Consider, } \|f(y + \beta a \mathbb{I}_d) - f(y)\|_{\infty}$$

by contraction property of f ,

$$\begin{aligned} \|f(y + a \mathbb{I}_d) - f(y)\|_{\infty} &\leq \beta \|y + a \mathbb{I}_d - y\|_{\infty} \\ &= \beta \|a \mathbb{I}_d\|_{\infty} \end{aligned}$$

$$\|\mathbb{I}_d\|_{\infty} = 1$$

$$\Rightarrow \|f(y + a \mathbb{I}_d) - f(y)\|_{\infty} \leq \beta |a|$$

$$\Rightarrow \max_i |(f(y + a \mathbb{I}_d) - f(y))(i)| \leq \beta |a|$$

(definition of ∞ norm)

$$\Rightarrow (f(y + a \mathbb{I}_d) - f(y))(i) \leq \beta |a| \quad \text{--- (1)}$$

$$\nexists (f(y + a \mathbb{I}_d) - f(y))(i) \geq -\beta |a| \quad \text{--- (2)}$$

$$\text{①} \Rightarrow f(y + a \mathbb{I}_d) \leq \beta \leq f(y) + \underbrace{\beta |a| \mathbb{I}_d}_{\text{--- (3)}}$$

$$\text{Note that } x \leq y + a \mathbb{I}_d$$

$$\Rightarrow f(x) \leq f(y + a \mathbb{I}_d) \quad (\because f(\cdot) \text{ is monotone}) \quad \text{--- (4)}$$

$$\text{③} \text{ and } \text{④} \Rightarrow f(x) \leq f(y + a \mathbb{I}_d) \leq f(y) + \beta |a| \mathbb{I}_d$$

$$\boxed{\Rightarrow f(x) \leq f(y) + \beta |a| \mathbb{I}_d.}$$

b)

\therefore Putting $a=0$, we get

if $x \leq y + \alpha \text{Id}$ then $f(x) \leq f(y) + \frac{\alpha}{\beta} \text{Id}$

$\Rightarrow x \leq y$ then $f(x) \leq f(y)$

\therefore if $x \leq y \Rightarrow f(x) \leq f(y)$ — ①

this shows that $f(\cdot)$ is monotone.

Consider 2 vectors $x \in Y$

we can always find a scalar a

such that $y - a \text{Id} \leq x \leq y + a \text{Id}$. — ②

$$(a \geq \max_i (y-x)(i), a \leq \max_i (y-x)(i))$$

$$\text{i.e. } a \geq \max_i |(y-x)(i)|$$

$$\Rightarrow a \geq \|y-x\|_{\infty}$$

so for finite $y \in X$ we can always find some scalar a .

Now apply

$$y - a \text{Id} \leq x$$

$$\Rightarrow y \leq x + a \text{Id}$$

Apply the given property,

$$f(y) \leq f(x) + |\beta| a \|Id.$$

$$\Rightarrow f(x) - f(y) \geq -|\beta| a \|Id \quad \text{--- (3)}$$

$$\text{Also } x \leq y + a \|Id \Rightarrow f(x) \leq f(y) + |\beta| a \|Id$$

$$\Rightarrow f(x) - f(y) \leq |\beta| a \|Id \quad \text{--- (4)}$$

$$(3) \& (4) \Rightarrow \|f(x) - f(y)\|_{\infty} \leq |\beta| a \|Id \quad \text{--- (5)}$$

$$(\because -(f(x) - f(y)) \leq |\beta| a \|Id \text{ (even)} \\ \& f(x) - f(y) \leq |\beta| a \|Id \text{ (odd)})$$

so all elements' absolute value $< |\beta| a \|Id$)

We showed that $a \geq \|x - y\|_{\infty}$

So putting $a = \|x - y\|_{\infty}$

$$(5) \Rightarrow \|f(x) - f(y)\|_{\infty} \leq |\beta| \|x - y\|_{\infty} \quad \text{--- (6)}$$

$f(\cdot)$ is a contraction wrt supremum norm

$f(\cdot)$ is a monotone contraction

To prove :

$$c) \quad x - \frac{1}{1-\beta} \|f(x) - x\|_{\infty} \leq f(x) - \frac{\beta}{1-\beta} \|f(x) - x\|_{\infty} I_d$$

— P1 (P1)

$$f(x) + \frac{\beta}{1-\beta} \|f(x) - x\|_{\infty} I_d \geq x + \frac{1}{1-\beta} \|f(x) - x\|_{\infty} I_d$$

— P2 (P2)

$$f(x) - \frac{\beta}{1-\beta} \|f(x) - x\|_{\infty} I_d \leq x^+ \leq f(x) + \frac{\beta}{1-\beta} \|f(x) - x\|_{\infty} I_d$$

— P3 (P3)

the above are the inequalities that have to be proved.

$x = f(x)$

$$\begin{aligned} x - f(x) &\leq \|x - f(x)\|_{\infty} I_d \\ x - f(x) &\geq -\|x - f(x)\|_{\infty} I_d \end{aligned} \quad \left. \begin{array}{l} \text{(definition} \\ \text{of max} \\ \text{norm)} \end{array} \right\}$$

— ① (1)
— ② (2)

$$\begin{aligned} ① \Rightarrow x - f(x) &\leq \left(\frac{1-\beta}{1-\beta} \|x - f(x)\|_{\infty} \right) I_d \\ &\quad \text{(multiply by } \frac{1}{1-\beta} \text{ divide by } \beta) \end{aligned}$$

$$\Rightarrow x \leq f(x) + \left(\frac{1}{1-\beta} - \frac{\beta}{1-\beta} \right) (\|x - f(x)\|_{\infty}) I_d.$$

$$\Rightarrow \underbrace{x - f(x)}_{\|x-f(x)\|_\infty \text{ Id}} \leq f(x) - \frac{\beta}{1-\beta} \|x-f(x)\|_\infty$$

$$\Rightarrow x - \frac{\|x-f(x)\|_\infty \text{ Id}}{1-\beta} \leq f(x) - \frac{\beta}{1-\beta} \|x-f(x)\|_\infty \quad \text{--- (3)}$$

P1 proved!

Similarly using (2) ,

$$x - f(x) \geq -\|x - f(x)\|_\infty \text{ Id} \left(\frac{1-\beta}{1-\beta} \right)$$

$$\Rightarrow x \geq f(x) - \left(\frac{1}{1-\beta} - \frac{\beta}{1-\beta} \|x-f(x)\|_\infty \text{ Id} \right)$$

$$\Rightarrow f(x) + \frac{\beta}{1-\beta} \|f(x)-x\|_\infty \text{ Id} \leq x + \frac{1}{1-\beta} \|f(x)-x\|_\infty \text{ Id} \quad \text{--- (4)}$$

P2 proved!

Consider $\|x - x^*\|_\infty$.

$$\begin{aligned} \|x - x^*\|_\infty &= \|x - f(x) - x^* + f(x^*)\|_\infty \\ &= \|-(f(x) - x) + (f(x) - x^*)\|_\infty \\ &\leq \|f(x) - x\|_\infty + \|f(x) - x^*\|_\infty \quad (\text{triangle inequality}) \\ \text{But } \|f(x) - f(x^*)\|_\infty &= \frac{\beta}{1-\beta} \|f(x) - f(x^*)\|_\infty \leq \beta \|x - x^*\|_\infty \end{aligned}$$

$\therefore x^*$ is fixed pt. & f(.) is unif.

$$\Rightarrow \|x - x^*\|_\infty \leq \|f(x) - x\|_\infty + \beta \|x - x^*\|_\infty$$

$$\therefore \|x - x^*\|_\infty \leq \frac{1}{1-\beta} \|f(x) - x\|_\infty$$

$$\text{So } \max_i |x - x^*| (i) \leq \frac{1}{1-\beta} \|f(x) - x\|_\infty$$

$$\Rightarrow x - x^* \leq + \frac{1}{1-\beta} \|f(x) - x\|_\infty \text{ Id} \quad \textcircled{A}$$

$$\& x - x^* \geq - \frac{1}{1-\beta} \|f(x) - x\|_\infty \text{ Id} \quad \textcircled{B}$$

$$\Rightarrow x \leq x^* + \frac{1}{1-\beta} \left(\frac{1}{1-\beta} \|f(x) - x\|_\infty \right) \text{ Id} \quad \textcircled{5}$$

$$\text{and } x^* \leq x + \frac{1}{1-\beta} \left(\frac{1}{1-\beta} \|f(x) - x\|_\infty \right) \text{ Id} \quad \textcircled{6}$$

these eqns are of the form $x \leq y + a \text{Id}$

$$\text{where } a = \frac{1}{1-\beta} \|f(x) - x\|_\infty$$

Applying the property from part (c), ($\because f$ is monotone contrac.)

$$\textcircled{8} \quad \Rightarrow f(x) \leq f(x^*) + \beta \left| \frac{1}{1-\beta} \|f(x) - x^*\|_\infty \right| \text{ Id.}$$

$$\text{But } f(x^*) = x^* \& \frac{1}{1-\beta} \|f(x) - x^*\|_\infty > 0$$

$$\Rightarrow f(x) \leq x^* + \frac{\beta}{1-\beta} \|f(x) - x^*\|_\infty \text{ Id}$$

Once again applying $x \leq y + a \text{ Id} \Rightarrow f(x) \leq f(y) + \frac{\beta}{1-\beta} \|f(y)-x\|_\infty \text{ Id}$

$$\textcircled{6} \Rightarrow f(x^+) \leq f(x) + \beta \left(\frac{1}{1-\beta} \|f(x)-x\|_\infty \text{ Id} \right) \text{ Id}$$

$$\therefore x^+ \leq f(x) + \frac{\beta}{1-\beta} \|f(x)-x\|_\infty \text{ Id} \quad \text{--- \textcircled{8}}$$

$$\textcircled{7} \Rightarrow x^+ \geq f(x) - \frac{\beta}{1-\beta} \|f(x)-x\|_\infty \text{ Id} \quad \text{--- \textcircled{9}}$$

$$\begin{aligned} \textcircled{8} \wedge \textcircled{9} \Rightarrow & f(x) - \frac{\beta}{1-\beta} \|f(x)-x\|_\infty \text{ Id} \leq x^+ \\ & \leq f(x) + \frac{\beta}{1-\beta} \|f(x)-x\|_\infty \text{ Id} \end{aligned} \quad \text{--- \textcircled{10}}$$

P₃ is proved!

$\exists P_1, \epsilon, P_2 \wedge P_3 \Rightarrow$

$$\begin{aligned} x - \frac{1}{1-\beta} \|f(x)-x\|_\infty \text{ Id} & \leq f(x) - \frac{\beta}{1-\beta} \|f(x)-x\|_\infty \text{ Id} \\ & \leq x^+ \leq f(x) + \frac{\beta}{1-\beta} \|f(x)-x^+\|_\infty \text{ Id} \\ & \leq x + \frac{1}{1-\beta} \|f(x)-x\|_\infty \text{ Id}. \end{aligned}$$

(eqns \textcircled{3}, \textcircled{4}, \textcircled{10})