

MDP components: state space  $\mathcal{X}$ , action space  $\mathcal{A}$ , transition probabilities  $p_{ij}(a)$

Expected cost or Value function

initial state

$$J_{\pi}(x_0) = \lim_{N \rightarrow \infty} E \left( \sum_{k=0}^{N-1} \alpha^k g(x_k, \mu_k(x_k), x_{k+1}) \right)$$

policy  $\pi = \{ \mu_0, \mu_1, \mu_2, \dots \}$

single-stage cost "stationary"

discount factor  $\alpha \in (0, 1]$

Goal:  $J^*(x_0) = \min_{\pi \in \Pi} J_{\pi}(x_0)$

optimal expected cost  $\rightarrow$  set of admissible policies

Let  $\pi^* = \arg \min_{\pi \in \Pi} J_{\pi}(x_0)$

optimal policy

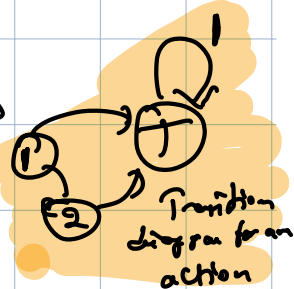
Two popular MDPs:

(I) Stochastic shortest path (SSP) aka episodic MDPs

- (a)  $\alpha = 1$
- (b) Finite state-action space  
states =  $\{1, \dots, n\} \cup \{T\}$

- (c) There exist a special state, say "T" that is
  - cost-free
  - absorbing
$$g(T, a, T) = 0, \forall a$$

$$P_{TT}(a) = 1, \forall a$$



## (II) Discounted MDP

(a)  $\alpha < 1$

(b)  $|g(i, a, j)| \leq M < \infty \quad \forall i, j \in \mathcal{X}, a \in \mathcal{A}$

(a) & (b)  $\Rightarrow J_{\pi}(x_0)$  is finite.

(III) Average-cost MDP : skipped. (see Chapter 5 of Vol II, DPoc book)

### Main results:

## (I) Taking finite horizon to the limit

Let  $J_N^*(i)$  be the optimal expected cost of a "N-stage" finite horizon MDP, with initial state  $i \in \mathcal{A}$  & stationary cost  $g(i, a, j)$

*Alert: this is a change of notation from finite horizon chapter*

Then, the infinite horizon optimal expected cost  $J^*$  is given by

$$J^*(i) = \lim_{N \rightarrow \infty} J_N^*(i)$$

$\uparrow$  infinite horizon optimal cost       $\uparrow$  finite horizon optimal cost

## (II) Bellman equation

Assume  $X = \{1, \dots, n\}$ , transition prob.  $P_{ij}(a)$

For a  $N$ -stage problem, with  $J_N^*$  denoting the optimal cost, the DP algorithm is

$$J_{k+1}^*(i) = \min_{a \in A(i)} \sum_{j=1}^n P_{ij}(a) (g(i, a, j) + \alpha J_k^*(j))$$

optimal cost in a  $(k+1)$ -stage problem

optimize action in the current stage

optimal cost for a  $k$ -stage problem

In the infinite horizon, the optimal cost  $J^*$  satisfies

$$J^*(i) = \min_{a \in A(i)} \sum_{j=1}^n P_{ij}(a) (g(i, a, j) + \alpha J^*(j))$$

(or)

Bellman equation

$$J^*(i) = \min_{a \in A(i)} E_{ij} [g(i, a, j) + \alpha J^*(j)], \forall i$$

DP algorithm is an algorithm, while Bellman equation is a system of eqns that the optimal cost satisfies.

(3) How to get the optimal policy  $\pi^*$ ?

$$\forall i, J^*(i) = \min_{a \in A(i)} \sum_{j=1}^n P_{ij}(a) (g(i, a, j) + \alpha J^*(j)) \leftarrow \text{Bellman equation}$$

For state  $i$ , Let  $a^*(i)$  be the minimizer in the Bellman equation

Make-up a policy as follows:  $\pi(i) = a^*(i), \forall i$

Then,  $\pi$  is ~~the~~ an optimal policy.

## Lecture-6 Stochastic shortest path (SSP) problems

State space =  $\{1, \dots, n\}$

In SSP,  $\exists$  a special terminal state, say  $T$ , that satisfies

$$P_{TT}(a) = 1, \quad g(T, a, T) = 0 \quad \forall a$$

absorbing                      cost-free

### Examples:

- ① A simple example of SSP: Deterministic shortest path
- ② A finite horizon problem can be easily cast as an SSP.

In the SSP for a finite horizon problem with horizon  $N$ :

States:  $(i, k)$

$\nearrow$  state       $\rightarrow$  stage

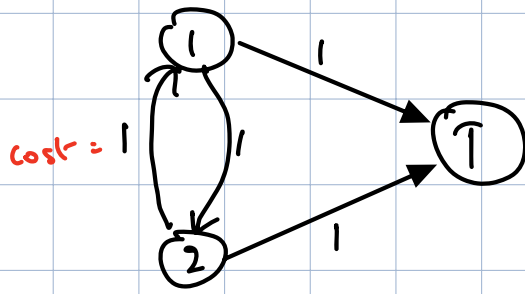
Terminal states:  $(i, N), \forall i \in X$

Transitions:  $(i, k) \xrightarrow{a} (j, k+1)$  w.p.  $P_{ij}(a)$

Costs:  $\pi.w.$



# Proper policies:



Shortest path problem

Improper policy: loop between 1 & 2

Expected cost =  $\infty$ .

Proper policy: Go to T from 1 & also 2.

**Stationary policy:**  $\pi: \mathcal{S} \rightarrow \mathcal{A}$  which takes the same action, say  $a$ , in a state  $i$ , irrespective of the stage  $k$  in the infinite horizon  $\pi = (\mu, \mu, \dots)$ . So, we identify the stationary policy with a mapping from  $\mathcal{S}$  to  $\mathcal{A}$ .

Def: A stationary policy  $\pi$  is **proper** if  $\exists m > 0$

$$P_{\pi} = \max_{i=1, \dots, n} P(\underbrace{i_m \neq T}_{\substack{\downarrow \\ \text{nth state} \\ \text{visited along} \\ \text{a sample path}}} \mid \underbrace{i_0 = i}_{\substack{\downarrow \\ \text{Starting in} \\ \text{state } i}}, \pi) < 1$$

$\downarrow$  actions governed by  $\pi$

Proper policy:  $\exists$  a pos. prob. path from any state  $i$  to the terminal state.

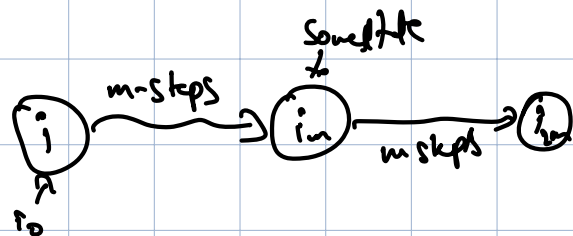
A policy that isn't proper is improper (i.e.,  $P_{\pi} = 1$ )

## Assumptions:

(A1) There exists at least one proper policy.

(A2) For every improper policy  $\pi$ , the associated expected cost  $J_{\pi}(i)$  is infinite for at least one state  $i$ .

Under a proper policy, Prob(not reaching T) goes down asymptotically.



$$P(i_{2m} \neq T \mid i_0 = i, \pi)$$

$$= P(i_{2m} \neq T \mid i_m \neq T, i_0 = i, \pi) P(i_m \neq T \mid i_0 = i, \pi)$$

Markov property

$$\rightarrow P(i_m \neq T \mid i_0 \neq T, \pi) P(i_m \neq T \mid i_0 = i, \pi) \leq e_{\pi}^2$$

More generally,

$$P(i_k \neq T \mid i_0 = i, \pi) \leq e_{\pi}^{\lfloor \frac{k}{m} \rfloor}$$

$$\begin{aligned} & \begin{matrix} \swarrow k < m & \searrow k > m \\ \leq & P(i_k \neq T \mid i_{\lfloor \frac{k}{m} \rfloor} \neq T, i_0 = i, \pi) \\ & \times P(i_{\lfloor \frac{k}{m} \rfloor} \neq T \mid i_0 = i, \pi) \\ & \leq 1 \times e_{\pi}^{\lfloor \frac{k}{m} \rfloor} = e_{\pi}^{\lfloor \frac{k}{m} \rfloor} \end{matrix} \end{aligned}$$

Let  $g_k$  be the cost incurred in  $k$ th stage for policy  $\pi$

$$E(g_k) \leq e_{\pi}^{\lfloor \frac{k}{m} \rfloor} \max_{i \in \mathcal{I}} |g(i, \pi(i))|$$

For simplicity, assume single stage cost doesn't depend on next state

$$|J_{\pi}(i)| \leq \lim_{N \rightarrow \infty} \sum_{k=0}^{N-1} e_{\pi}^{\lfloor \frac{k}{m} \rfloor} \max(g(i, \pi(i))) < \infty$$

So, for a proper policy  $\pi$ , the expected cost  $J_{\pi}$  is finite.

## Lecture-7\*

### Notation:

State space  $\mathcal{X} = \{1, \dots, n, T\}$

For any  $J = (J(1), \dots, J(n))$ ,  
define the "Bellman optimality" operator  $TJ = (TJ(1), \dots, TJ(n))$

as

For  $i=1, \dots, n$ ,

$$(TJ)(i) = \min_{a \in \mathcal{A}(i)} \sum_{j \in \{1, \dots, n, T\}} p_{ij}(a) g(i, a, j) + \sum_{j=1}^n p_{ij}(a) J(j)$$

operator & it is not tied to any policy

includes terminal state

does not include "T"

Equivalently,

$$(TJ)(i) = \min_{a \in \mathcal{A}(i)} E_j (g(i, a, j) + J(j))$$

Another operator for a given policy  $\pi$ :

$T_{\pi} J = (T_{\pi} J(1), \dots, T_{\pi} J(n))$ , where

$$T_{\pi} J(i) = \sum_{j \in \mathcal{X}} p_{ij}(\pi(i)) (g(i, \pi(i), j) + J(j)), \quad i=1, \dots, n$$

For a given  $\pi$ ,

$$P^\pi = \begin{bmatrix} P_{11}(\pi(1)) & \dots & P_{1n}(\pi(1)) \\ \vdots & & \vdots \\ P_{n1}(\pi(n)) & \dots & P_{nn}(\pi(n)) \end{bmatrix}$$

$P^\pi$  matrix has positive entries & row sum  $\leq 1$

$$g_\pi = (g_\pi(1), \dots, g_\pi(n)), \text{ where}$$

$$g_\pi(i) = \sum_{j \in X} P_{ij}(\pi(i)) g(i, \pi(i), j)$$

Using  $g_\pi$  &  $P^\pi$ , we have

$$T_\pi J = g_\pi + P^\pi J$$

Applying operator  $T$  multiple times:

$$J \xrightarrow{T} TJ \xrightarrow{T} T^2 J \rightarrow \dots$$

$$\text{Let } (T^k J)(i) = (T(T^{k-1} J))(i), \quad i=1 \dots n$$

with  $(T^0 J)(i) = J(i)$

Similarly for  $T_\pi$ .

" $T$  &  $T_\pi$  are monotone"

Lemma 1: Let  $J, J' \in \mathbb{R}^n$  & satisfy  $J(i) \leq J'(i), \forall i$

Then, for any  $k=1, 2, \dots$

$$(i) (T^k J)(i) \leq (T^k J')(i), \text{ and}$$

(ii) For any stationary policy  $\pi$ ,

$$(T_\pi^k J)(i) \leq (T_\pi^k J')(i)$$

Pf:

$$\begin{aligned} (TJ)(i) &= \min_a \sum_{j \in X} P_{ij}(a) (g(i, a, j) + J(j)) \\ &\leq \min_a \sum_{j \in X} P_{ij}(a) (g(i, a, j) + J'(j)) \\ &= (TJ')(i). \end{aligned}$$

**H.W.** Complete the rest of the proof using induction.

**H.W.** Do the proof for  $T_\pi$ . ■

Another lemma: (Constant-shift lemma)

Stationary  $\pi$ ,  $\delta \rightarrow$  any scalar,  $e \rightarrow$  vector of  $n$  ones.

Then,  $\forall i=1, \dots, n$ ,  $\forall k=1, 2, \dots$ , we have

$$(i) (T^k (J + \delta e))(i) \leq (T^k J)(i) + \delta$$

$$(ii) (T_\pi^k (J + \delta e))(i) \leq (T_\pi^k J)(i) + \delta.$$

$$J + \delta e = \begin{bmatrix} J(1) \\ \vdots \\ J(n) \end{bmatrix} + \begin{bmatrix} \delta \\ \vdots \\ \delta \end{bmatrix}$$

Pf:

$$\begin{aligned}
 & (T(\mathcal{J} + \delta e))(i) \\
 &= \min_a \sum_{j \in X} p_{ij}(a) (g(i, a, j) + \sum_{\tilde{j}=1}^n p_{\tilde{j}j}^{(a)} (\mathcal{J} + \delta e)(\tilde{j})) \\
 &= \min_a \sum_{j \in X} p_{ij}(a) g(i, a, j) + \sum_{\tilde{j}=1}^n p_{\tilde{j}j}^{(a)} (\mathcal{J}(\tilde{j}) + \delta) \\
 &= \min_a \left[ \sum_j p_{ij}(a) (g(i, a, j) + \mathcal{J}(j)) + \sum_{\tilde{j}=1}^n p_{\tilde{j}j}(a) \delta \right] \leq 1
 \end{aligned}$$

Since  $\sum_{\tilde{j}=1}^n p_{\tilde{j}j}(a) \leq 1$   $\rightarrow$   $\leq 1$

no  $\mathcal{J}(i)$  here

$$\begin{aligned}
 & \leq \left[ \min_a \sum_{j \in X} p_{ij}(a) g(i, a, j) + \sum_{\tilde{j}=1}^n p_{\tilde{j}j}(a) \mathcal{J}(\tilde{j}) \right] + \delta \\
 &= (T\mathcal{J})(i) + \delta
 \end{aligned}$$

We showed

$$(T(\mathcal{J} + \delta e))(i) = (T\mathcal{J})(i) + \delta$$

Rest of the proof  $\rightarrow$  a simple induction.  $\Leftarrow$  (H.W.)

### Properties of $T_\pi$ :

Proposition 1: Assume (A1) & (A2). Then,

$\rightarrow$  No proper policy

$\rightarrow$  improper policies have infinite cost

(i) For any proper policy  $\pi$ , the associated cost  $J_\pi$  satisfies

(\*\*)  $\rightarrow \lim_{k \rightarrow \infty} (T_\pi^k J)(i) = J_\pi(i), i=1, \dots, n,$

for any  $J$ .

Also,  $J_\pi = T_\pi J_\pi$  — (\*\*)

&  $J_\pi$  is the unique solution of (\*\*)

$$J(i) = g^{i, \pi(i)} + \sum p_{ij}(\pi(i)) J(j), \forall i$$

Policy evaluation: Task of computing  $J_\pi$  for a given policy  $\pi$

Using (\*\*), policy evaluation can be done by starting in a "J" & repeatedly applying  $T_\pi$  ( $J \xrightarrow{T_\pi} T_\pi J \xrightarrow{T_\pi} T_\pi^2 J \rightarrow \dots$ )

Using (\*\*), solve " $J = T_\pi J$ " to obtain  $J_\pi$ .

(ii) Suppose a stationary policy  $\pi$  satisfies

$$J(i) \geq (T_\pi J)(i), \quad i=1 \dots n, \quad \text{for some finite } J.$$

Then,  $\pi$  is proper.

Properties of Bellman optimality operator  $T$ :

Prop 2: Assume (A1) & (A2).

$$\text{Let } J^*(i) = \min_{\pi \in \Pi} J_\pi(i), \quad \forall i=1 \dots n$$

$\pi^* \in$  optimal policy

Then,

(i)  $J^*$  satisfies

Bellman equation  $\longrightarrow$   $J^* = T J^*$  — (\*\*\*)

&  $J^*$  is the unique solution to (\*\*\*)

i.e., does not exist a  $\tilde{J} \neq J^*$  satisfying  $\tilde{J} = T \tilde{J}$

$$f(x) = x/2$$

$$f(x^*) = x^*$$

$$16 \xrightarrow{f} 8 \xrightarrow{f} 4 \xrightarrow{f} 2 \xrightarrow{f} \dots \rightarrow 0$$

(ii) For any finite  $J$ , we have

$$\lim_{k \rightarrow \infty} (T^k J)(i) = J^*(i), \quad i=1, \dots, n$$

(iii) A stationary policy  $\pi$  is optimal if & only if

$$T_\pi J^* = T J^*$$

$$\sum_j P_{ij}(\pi(i)) (g(i, \pi(i), j) + J^*(j)) = \min_a \sum_j P_{ij}(a) (g(i, a, j) + J^*(j)), \quad \forall i$$

Lecture-8\*

< Proof of Prop. 1 > part (i)

First claim:  $\lim_{k \rightarrow \infty} T_\pi^k J = J_\pi$  for any  $J$

< cf. within pf >  $T_\pi J = g_\pi + P_\pi J$

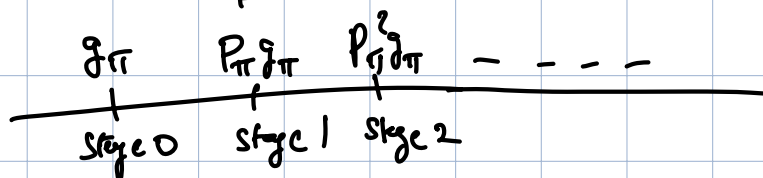
$$T_\pi^2 J = g_\pi + P_\pi T_\pi J = g_\pi + P_\pi g_\pi + P_\pi^2 J$$

Generalizing, we obtain

$$T_\pi^k J = P_\pi^k J + \sum_{m=0}^{k-1} P_\pi^m g_\pi \quad \text{--- (†)}$$

Note:  $\lim_{k \rightarrow \infty} \sum_{m=0}^{k-1} P_\pi^m g_\pi = J_\pi$  --- (‡)

(Recall  $J_\pi(i) = \lim_{k \rightarrow \infty} E \left( \sum_{m=0}^{k-1} g(x_m, \pi(x_m), x_{m+1}) \mid x_0 = i \right)$ )



← infinite horizon  
sum of these  $P_\pi^m g_\pi$   
is  $J_\pi$



$\pi$  is proper  $P(i_k \neq T | i_0 = i, \pi) \leq e^{-\frac{k}{\pi}}$ ,  $\forall i$

$$\lim_{k \rightarrow \infty} P_{\pi}^k J = 0 \quad \text{--- (I)}$$

Using (I) & (II) in (+), we obtain  $\lim_{k \rightarrow \infty} T_{\pi}^k J = J_{\pi}$

Second claim:  $J_{\pi} = T_{\pi} J_{\pi}$  & uniqueness.

$$T_{\pi}^{k+1} J = g_{\pi} + P_{\pi} T_{\pi}^k J$$

Taking limit as  $k \rightarrow \infty$  on both sides, we get

$$\lim_{k \rightarrow \infty} T_{\pi}^{k+1} J = \lim_{k \rightarrow \infty} (g_{\pi} + P_{\pi} T_{\pi}^k J)$$

$$J_{\pi} = g_{\pi} + P_{\pi} J_{\pi} = T_{\pi} J_{\pi}$$

$$\text{So, } J_{\pi} = T_{\pi} J_{\pi}$$

Uniqueness: Suppose  $\tilde{J}$  satisfies  $\tilde{J} = T_{\pi} \tilde{J}$

$$\tilde{J} = T_{\pi} \tilde{J} = T_{\pi}^2 \tilde{J} = \dots = T_{\pi}^k \tilde{J}$$

Taking limits,  $\tilde{J} = \lim_{k \rightarrow \infty} T_{\pi}^k \tilde{J} = J_{\pi}$

So,  $J_{\pi}$  is the unique fixed point.

< End of Pt within Pt, again >

< Proof of Prop 1 - part (ii) >

Given: Stationary  $\pi$  satisfying  $J \geq T_{\pi} J$  for some  $J$  finite

$$J \geq T_{\pi} J \Rightarrow T_{\pi} J \geq T_{\pi}^2 J$$

$$\Rightarrow J \geq T_{\pi}^2 J \quad \& \quad \text{so on}$$

$$J \geq T_{\pi}^k J = P_{\pi}^k J + \sum_{m=0}^{k-1} P_{\pi}^m g_{\pi} \quad \left( \begin{array}{l} \text{From eq (4)} \\ \text{in part (i)} \\ \text{proof} \end{array} \right)$$

If  $\pi$  is not proper, then

$$J_{\pi} = \lim_{k \rightarrow \infty} \sum_{m=0}^{k-1} P_{\pi}^m g_{\pi} \quad \text{diverges}$$

This leads to a contradiction since  $J \geq \lim_{k \rightarrow \infty} T_{\pi}^k J = J_{\pi}$

$\uparrow$   
finite
 $\uparrow$   
infinit for unproper policy

< End of part (ii) >

< Proof of Prop 2 >

(i) <sup>to show:</sup>  $J^*$  satisfies

$$J^* = T J^* \quad (**)$$

&  $J^*$  is the unique solution to (\*\*)

< Pf > "T has at most one fixed point."

Suppose  $J$  and  $J'$  satisfy  $J = T J$ ,  $J' = T J'$

We will show  $J \geq J'$  by an interchangeable argument which would imply  $J = J'$

Select a policy  $\pi$  s.t.  $T_{\pi} J = T J$

$$\sum_j P_{ij}(\pi(i)) (g(i, \pi(i), j) + J(j)) = \min_a \sum_j P_{ij}(a) (g(i, a, j) + J(j))$$

"Set  $\pi(i)$  to be the minimizing action on the RHS"

Similarly, make up a policy  $\pi'$  s.t.  $T J' = T_{\pi'} J'$

We have  $T = T \mathcal{J} = T_{\pi} \mathcal{J}$

If  $T = T_{\pi} \mathcal{J}$ , then  $\mathcal{J} = \mathcal{J}_{\pi}$ . Also,  $\pi$  is proper  
(any part (ii) after observing  $\mathcal{J} = T_{\pi} \mathcal{J}$ )

Similarly,  $\mathcal{J}' = \mathcal{J}_{\pi'}$  &  $\pi'$  is proper.

Now,  $T = T \mathcal{J} = T^k \mathcal{J} \leq T_{\pi'}^k \mathcal{J}$ ,  $\forall k \geq 1$

(compare  $T \mathcal{J} \leq T_{\pi'} \mathcal{J}$  for any  $\pi'$ )  
 $\min_a \sum_i p_{ij}(a) (g(i,a,j) + \mathcal{J}(j)) \leq \sum_i p_{ij}(\pi'(i)) (g(i, \pi'(i)) + \mathcal{J}(j))$

$\mathcal{J} \leq T_{\pi'}^k \mathcal{J} \xrightarrow[\text{limits}]{\text{Take}}$   $\mathcal{J} \leq \lim_{k \rightarrow \infty} T_{\pi'}^k \mathcal{J} = \mathcal{J}_{\pi'} = \mathcal{J}'$

We got  $\mathcal{J} \leq \mathcal{J}'$

Repeating the arguments with  $\mathcal{J}$  &  $\mathcal{J}'$  swapped leads to  $\mathcal{J}' \leq \mathcal{J}$

So,  $\mathcal{J} = \mathcal{J}'$  (or)  $T$  has at most one fixed point.

< End of "uniqueness" part >

< Begin of "existence" part >

Let  $\pi$  be a proper policy (A1 ensures  $\exists$  at least one proper policy)

Let  $\pi'$  be another policy s.t.

$$T_{\pi'} \mathcal{J}_{\pi} = T \mathcal{J}_{\pi}$$

$$\mathcal{J}^* + \delta e = T(\mathcal{J}^* + \delta e) = \lim_{k \rightarrow \infty} T^k(\mathcal{J}^* + \delta e) = \mathcal{J}^* \Rightarrow \delta = 0$$

Now,  $J_\pi \stackrel{\text{Prop 1}}{=} T_\pi J_\pi \stackrel{\text{def of } T, J_\pi}{\geq} T J_\pi \stackrel{\text{by construction}}{=} T_{\pi'} J_\pi$

So,  $J_\pi \geq T_{\pi'} J_\pi \geq T_{\pi'}^2 J_\pi \dots \geq T_{\pi'}^k J_\pi$

Taking limits,  $J_\pi \geq \lim_{k \rightarrow \infty} T_{\pi'}^k J_\pi = J_{\pi'}$

Also,  $\pi'$  is proper from Prop 1, part (ii) since  $J_\pi \geq T_{\pi'} J_\pi$  &  $\pi$  is proper.

So, we have two proper policies satisfy

$$J_\pi \geq J_{\pi'}$$

Repeating the arguments above again & again & again, we obtain a sequence of policies  $\{\pi_k\}$

s.t. (i) Each  $\pi_k$  is proper

(ii)  $J_{\pi_k} \geq T J_{\pi_k} \geq J_{\pi_{k+1}}, \forall k \geq 1$

No. of policies is finite, because we assumed that the state & action spaces are finite.

$$J_{\pi_1} \geq T J_{\pi_1} \geq J_{\pi_2} \dots$$

in this sequence, a policy has to repeat

For that repeat policy, say  $\tilde{\pi}$ , we get

$$J_{\tilde{\pi}} = T J_{\tilde{\pi}}$$

$\Rightarrow$  fixed point exists.

< End of existence proof >

Lecture 9 Claim:  $T^k J \rightarrow J_{\tilde{\pi}}$  as  $k \rightarrow \infty$  (later we show  $J_{\tilde{\pi}} = J^*$ )

<pf>:  $e = n$ -vector of ones. Fix  $\delta > 0$   
 Pick  $\hat{J}$  s.t.  $T_{\tilde{\pi}} \hat{J} = \hat{J} - \delta e$

Such a  $\hat{J}$  can be found since

$$\hat{J} = T_{\tilde{\pi}} \hat{J} + \delta e = (g_{\tilde{\pi}} + \delta e) + P_{\tilde{\pi}} \hat{J}$$

$$\hat{J} = (g_{\tilde{\pi}} + \delta e) + P_{\tilde{\pi}} \hat{J} \quad (*)$$

$$J_{\tilde{\pi}} = g_{\tilde{\pi}} + P_{\tilde{\pi}} J_{\tilde{\pi}} = T_{\tilde{\pi}} J_{\tilde{\pi}}$$

$\hat{J}$  is the expected cost of  $\tilde{\pi}$  in an SSP with single stage cost  $(g_{\tilde{\pi}} + \delta e)$

$J_{\tilde{\pi}}$  is expected cost in an SSP with single stage cost  $g_{\tilde{\pi}}$ . This is sol for SSP with  $g_{\tilde{\pi}}$

$\hat{J}$  is the unique soln to (\*) since  $\tilde{\pi}$  is proper.

By construction,  $J_{\tilde{\pi}} \leq \hat{J}$  since  $g_{\tilde{\pi}} \leq g_{\tilde{\pi}} + \delta e$

Recall  $J_{\tilde{\pi}}$  satisfies  $J_{\tilde{\pi}} = T J_{\tilde{\pi}}$  ← from previous proof

Notice that

$$J_{\tilde{\pi}} = T J_{\tilde{\pi}} \leq T \hat{J} \leq T_{\tilde{\pi}} \hat{J} = \hat{J} - \delta e \leq \hat{J}$$

$\uparrow$   $T$  is monotone  
 $\uparrow$   $T$  is Bellman optimality operator  $\Rightarrow T \leq T_{\tilde{\pi}} \forall \pi$ .  
 $\uparrow$  trivial  
 $\uparrow$  By construction of  $\hat{J}$

$$\Rightarrow J_{\tilde{\pi}} = T^k J_{\tilde{\pi}} \leq T^k \hat{J} \leq T^{k-1} \hat{J} \leq \hat{J}$$

repeated application of  $T$  in earlier inequality

$$\begin{aligned} T \hat{J} &\leq \hat{J} \\ T^{k-1} (T \hat{J}) &\leq T^{k-1} \hat{J} \\ T^k \hat{J} &\leq T^{k-1} \hat{J} \end{aligned}$$

$J_1 \leq J_2$   
 $\Rightarrow T J_1 \leq T J_2$   
 contraction lemma

$$\Rightarrow J_{\text{ff}} \leq T^k \hat{J} \leq \hat{J} \quad \text{are both finite}$$

So,  $\{T_k \hat{J}\}$  forms a monotone bounded sequence.

Hence,  $\lim_{k \rightarrow \infty} T^k \hat{J} = \tilde{J}$ , for some  $\tilde{J}$

Apply  $T$  on both sides of the shaded equation

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

$T$  is a continuous mapping  $\Rightarrow T$  can be taken inside the limit

$TJ = \min_a \text{"linear-fg"} \Rightarrow T$  is continuous

$$T \tilde{J} = T \left( \lim_{k \rightarrow \infty} T^k \hat{J} \right) = \lim_{k \rightarrow \infty} T^{k+1} \hat{J} = \tilde{J}$$

So,  $T \tilde{J} = \tilde{J}$  (\*)

$\Rightarrow \tilde{J} = J_{\text{ff}}$  since the fixed point of  $T$  is unique &  $= J_{\text{ff}}$ .

"The sandwich principle"

We will show  $\lim_{k \rightarrow \infty} T^k (J_{\text{ff}} - \delta \epsilon) = J_{\text{ff}}$

To see this,

See a lemma on p. 9 above

$$J_{\pi} - \delta c = T J_{\pi} - \delta c \leq T (J_{\pi} - \delta c) \leq T J_{\pi} = J_{\pi}$$

$J_{\pi} = T J_{\pi}$  (pointing to the first  $J_{\pi}$ )  
 $T$  is monotone  
 $J_{\pi} - \delta c \leq J_{\pi}$  (pointing to the inequality)

$$J_{\pi} - \delta c \leq T (J_{\pi} - \delta c) \leq J_{\pi}$$

Using "  $T(J_{\pi} - \delta c) \leq T^2(J_{\pi} - \delta c) \leq \dots \leq T^k(J_{\pi} - \delta c)$  "

$$J_{\pi} - \delta c \leq T^k(J_{\pi} - \delta c) \leq J_{\pi}$$

So,  $\{T^k(J_{\pi} - \delta c)\}$  is a monotone bounded sequence, implying

$$\lim_{k \rightarrow \infty} T^k(J_{\pi} - \delta c) = J_{\pi}$$

Use an argument similar to (\*)

For any  $J$ , we can find a  $\delta > 0$  such that

$$J_{\pi} - \delta c \leq J \leq \hat{J}$$

cost of policy  $\pi$  with single stage cost  $g_{\pi} + \delta c$

$$T^k(J_{\pi} - \delta c) \leq T^k J \leq T^k \hat{J}$$

Taking limits,

$$\lim_{k \rightarrow \infty} T^k(J_{\pi} - \delta c) \leq \lim_{k \rightarrow \infty} T^k J \leq \lim_{k \rightarrow \infty} T^k \hat{J}$$

$$J_{\pi} \leq \lim_{k \rightarrow \infty} T^k J \leq J_{\pi} \Rightarrow \lim_{k \rightarrow \infty} T^k J = J_{\pi}$$

for any finite  $J$

What remains: To show  $J_{\pi} = J^*$

Take any policy  $\pi'$ . We have

$$T_{\pi'}^k J_0 \geq T^k J_0, \quad (*)$$

where  $J_0$  is an arbitrary  $n$ -vector.

Take limits as  $k \rightarrow \infty$  on both sides of  $(*)$  to obtain

$$J_{\pi'} \geq J_{\pi} \quad \text{--- True for any } \pi'$$

$$\text{So, } J_{\pi} = J^*.$$

< End of part (b), we showed (i)  $J^* = T J^*$   
(ii)  $\lim_{k \rightarrow \infty} T^k J = J^*$   
for any  $J$  >

< Start of part (c):

$$\pi \text{ is optimal } \Leftrightarrow T_{\pi} J^* = T J^* >$$

( $\Rightarrow$ ) Suppose  $\pi$  is optimal

$$J_{\pi} = J^*$$

$$T_{\pi} J^* = T_{\pi} J_{\pi} = J_{\pi} = J^* = T J^*$$

$$\text{So, } T_{\pi} J^* = T J^*$$



( $\Leftarrow$ ) Suppose  $J^* = T J^* = T_\pi J^*$

$J^* = T_\pi J^* \Rightarrow$  From Prop 1 part (ii),  $\pi$  is proper

For a proper  $\pi$ , we have

$$J^* = T_\pi J^*$$

$$\Rightarrow J^* = J_\pi$$

$\Rightarrow \pi$  is optimal.

< End of part (c) >

## Lecture-10\*

Example: Time to termination

Consider an SSP where

$$g(i, a) = 1 \quad i=1 \dots n$$

Goal: get to terminal state asap.

Let  $J^* \rightarrow$  optimal expected cost satisfies

Bellman equation:

$$J^*(i) = T J^*(i), \quad \forall i$$

$$J^*(i) = \min_a \left[ 1 + \sum_{j=1}^n P_{ij}(a) J^*(j) \right], \quad i=1 \dots n$$

Special case: Only one action in each state ← Markov chain

$J^*(i)$  ← expected first passage time to terminal state "T".

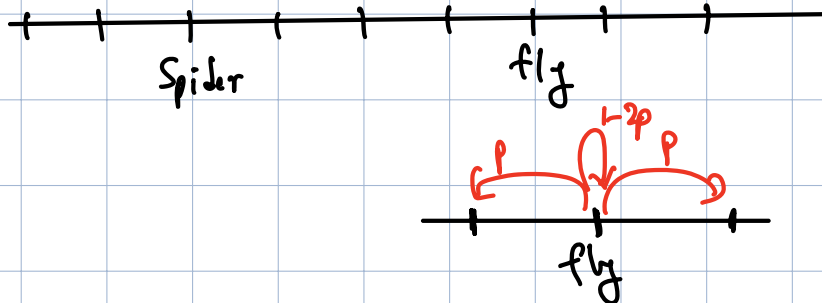
Let  $m_i$  denote this time (instead of  $J^*(i)$ )

$$m_i = 1 + \sum_{j=1}^n P_{ij} m_j, \quad i=1, \dots, n$$



Another example: Fly & a spider

Spider & fly move on a line.



Spider: (i) if at a distance  $> 1$  from fly,  
then jump 1 unit towards fly.

(ii) if distance  $\geq 1$

↙ jump towards fly  
↘ don't jump.

If fly & spider in same position, then you know what happens.

Goal: help spider catch the fly within "min" expected time.

SSP formulation!

State = distance between spider & fly

States =  $\{0, 1, \dots, n\}$  ↪ initial distance b/w spider & fly

Terminal state = 0

let  $P_{ij}(M)$  &  $P_{ij}(\bar{M})$

↑  
spider moves

↓  
spider doesn't move

Cost: = 1 in all states before hitting terminal "0" state.

remaining transition probabilities:  $P_{ij}$   $i \geq 2$

For  $i \geq 2$ ,  $P_{i,i} = p$ ,  $P_{i,i-1} = 1-2p$ ,  $P_{i,i-2} = p$

For  $i=1$ , "Spider jumps" action  $M$

$$P_{11}(M) = 2p$$

$$P_{10}(M) = 1-2p$$

"Spider doesn't jump" action  $\bar{M}$

$$P_{12}(\bar{M}) = p,$$

$$P_{11}(\bar{M}) = 1-2p, \quad P_{10}(\bar{M}) = p$$

Bellman equation for  $i \geq 2$

$$J^*(i) = 1 + p J^*(i) + (1-2p) J^*(i-1) + p J^*(i-2)$$

$$J^*(0) = 0$$

(\*)

For state = 1 :

$$J^*(1) = 1 + \min \left[ \begin{array}{l} 2p J^*(1) + (1-2p) J^*(0), \\ p J^*(2) + (1-2p) J^*(1) + p J^*(0) \end{array} \right]$$

$$J^*(1) = 1 + \min \left[ \underbrace{2p J^*(1)}_{\text{jump}}, \underbrace{p J^*(2) + (1-2p) J^*(1)}_{\text{don't jump}} \right] \quad (*)$$

From (\*),

$$J^*(2) = 1 + p J^*(2) + (1-2p) J^*(1)$$

$$\Rightarrow J^*(2) = \frac{1}{1-p} + \frac{(1-2p) J^*(1)}{1-p} \quad (**)$$

Substitute (\*\*) in (\*),

$$J^*(1) = 1 + \min \left[ 2p J^*(1), \frac{p}{1-p} + \frac{p(1-2p) J^*(1)}{1-p} + (1-2p) J^*(1) \right]$$

which is the same as

$$J^*(1) = 1 + \min \left[ \underbrace{2p J^*(1)}_{(A) \text{ jump}}, \underbrace{\frac{p}{1-p} + \frac{(1-2p) J^*(1)}{1-p}}_{(B) \text{ don't jump}} \right]$$

Case (A)  $\leq$  (B):

$$J^*(1) = 1 + 2p J^*(1) \quad \&$$

$$2p J^*(1) \leq \frac{p}{1-p} + \frac{(1-2p) J^*(1)}{1-p} \quad \text{--- (****)}$$

In this case,  $J^*(1) = \frac{1}{1-2p}$

Substitute this in (\*\*\*\*) to obtain

$$\frac{2p}{1-2p} \leq \frac{p}{1-p} + \frac{1}{1-p} \quad (\Leftrightarrow) \quad p \leq \frac{1}{3} \quad \text{--- optimal to jump.}$$

Case (A) > (B)

$$J^*(1) = 1 + \frac{p}{1-p} + \frac{(1-2p) J^*(1)}{1-p}, \quad \text{and (****)}$$

$$2p J^*(1) > \frac{p}{1-p} + \frac{(1-2p) J^*(1)}{1-p}$$

From (\*\*\*\*),  $J^*(1) = \frac{1}{p}$

Substitute this in the constraint to obtain

$$2 > \frac{p}{1-p} + \frac{1-2p}{p(1-p)} \quad (\Leftrightarrow) \quad p > \frac{1}{3}$$

↑  
optimal to not jump.

$$J^*(1) = \begin{cases} \frac{1}{1-2p} & \text{when } p \leq \frac{1}{3} \\ \frac{1}{p} & \text{else} \end{cases}$$

Using  $J^*(1)$ , we obtain  $J^*(2)$  & so on for  $J^*(i)$ ,  $i \geq 3$ .

Lecture 11

Coming next!

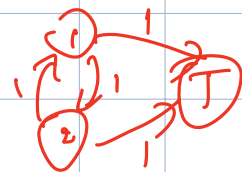
### Value iteration (VI)

Start with  $J_0$

$$J_0 \xrightarrow{T} J_1 \xrightarrow{T} J_2 \rightarrow \dots$$

$$J_{k+1} = T J_k$$

$$\begin{aligned} J_1 &= T J_0 \\ J_2 &= T J_1 \\ &= T^2 J_0 \\ &\dots \end{aligned}$$



We already know very clearly that

$$\lim_{k \rightarrow \infty} T^k J_0 = J^*$$

under (A1) & (A2)

$\exists$  proper policy

improper policies have  $\infty$  cost

"Can do VI for policy evaluation"

$$J_0 \xrightarrow{T_\pi} J_1 \xrightarrow{T_\pi} J_2 \xrightarrow{T_\pi} \dots \xrightarrow[k \rightarrow \infty]{} J_\pi$$

Special case: **All policies are proper**

Then, VI converges at a geometric rate.

$$\|J_k - J^*\|_\xi \leq e^{-k} \|J_0 - J^*\|_\xi$$

$\rightarrow$  iterations of VI

$\leftarrow$  need theory of contraction mappings.

mappings.

$\|\cdot\|_\xi \rightarrow$  weighted-max norm

Note! Can find a  $k$  s.t.  $\|J_k - J^*\| \leq \epsilon$ , for any  $\epsilon > 0$

$0 < c < 1$

constant

## Lecture-11\*

"Assume all policies are proper".

We will show that the Bellman optimality operator  $T$  is a contraction w.r.t. a weighted max-norm.

In particular,  $\exists$  a vector  $\xi = (\xi(1), \dots, \xi(n))$  s.t.  $\xi(i) > 0 \forall i=1, \dots, n$ , and a scalar  $0 < \rho < 1$  such that

$$\|TJ - T\bar{J}\|_{\xi} \leq \rho \|J - \bar{J}\|_{\xi}, \quad \forall J, \bar{J} \in \mathbb{R}^n$$

modulus of contraction  
↓

Here

$$\|J\|_{\xi} = \max_{i=1, \dots, n} \frac{|J(i)|}{\xi(i)}$$

↓  
weighted max-norm

same  $\rho$  for  $T$  &  $T_{\pi}$

Also, for any stationary proper policy  $\pi$ ,

$$\|T_{\pi}J - T_{\pi}\bar{J}\|_{\xi} \leq \rho \|J - \bar{J}\|_{\xi} \quad \forall J, \bar{J} \in \mathbb{R}^n$$

PF:

Need: a vector  $\xi$  such that  $T$  is a contraction w.r.t.  $\|\cdot\|_{\xi}$

Idea: Make up a MDP, solve it & get  $\xi$ .

& then show such a  $\xi$  helps in contraction business.

Consider a new SSP with same " $P_{ij}(a)$ ",  $\forall i, j=$   
but, different transition costs.

Transition cost = -1 everywhere except terminal state

$$g(T, a, T) = 0 \quad \forall a \quad \& \quad g(i, a, j) = -1 \quad \text{else}$$

Let  $\hat{J} \rightarrow$  optimal cost in this new SSP.  
 Then, using Bellman equation

$$\hat{J}(i) = -1 + \min_a \sum_{j \in X} P_{ij}(a) \hat{J}(j)$$

$$= T \hat{J}(i)$$

$$\leq T_{\pi} \hat{J}(i) \quad \text{for any proper } \pi$$

$$= -1 + \sum_j P_{ij}(\pi(i)) \hat{J}(j)$$

Let  $\xi(i) = -\hat{J}(i)$ ,  $i = 1 \dots n$

Note: ①  $\xi(i) \geq 1$

②  $-\hat{J}(i) \geq -1 + \sum_j P_{ij}(\pi(i)) (-\hat{J}(j))$

$$\xi(i) \geq 1 + \sum_j P_{ij}(\pi(i)) \xi(j) \quad (*)$$

$$\sum_j P_{ij}(\pi(i)) \xi(j) \leq \xi(i) - 1 \leq \rho \xi(i), \quad (**)$$

follows from (\*)

where

$$\rho = \max_{i=1 \dots n} \left( \frac{\xi(i) - 1}{\xi(i)} \right) < 1$$



Now, for any  $\pi, \mathcal{J}, \bar{\mathcal{J}}$ ,

$$T_{\pi} \mathcal{J}(i) = \sum_j P_{ij}(\pi(i)) (\mathcal{J}(j) + \tau(i, j))$$

$$|(T_{\pi} \mathcal{J})(i) - (T_{\pi} \bar{\mathcal{J}})(i)| = \left| \sum_{j=1}^n P_{ij}(\pi(i)) (\mathcal{J}(j) - \bar{\mathcal{J}}(j)) \right|$$

triangle inequality

$$\leq \sum_{j=1}^n P_{ij}(\pi(i)) |\mathcal{J}(j) - \bar{\mathcal{J}}(j)|$$

$$= \sum_j P_{ij}(\pi(i)) \xi(j) \left( \frac{|\mathcal{J}(j) - \bar{\mathcal{J}}(j)|}{\xi(j)} \right)$$

$$\begin{aligned} |a-b| &\leq 10 \\ a &\leq b+10 \\ b &\leq a+10 \end{aligned}$$

$$\leq \left( \sum_j P_{ij}(\pi(i)) \xi(j) \right) \left( \max_{j=1 \dots n} \frac{|\mathcal{J}(j) - \bar{\mathcal{J}}(j)|}{\xi(j)} \right)$$

$$|(T_{\pi} \mathcal{J})(i) - (T_{\pi} \bar{\mathcal{J}})(i)| \leq e \xi(i) \max_{j=1 \dots n} \frac{|\mathcal{J}(j) - \bar{\mathcal{J}}(j)|}{\xi(j)}$$

(I)

$$\frac{|(T_{\pi} \mathcal{J})(i) - (T_{\pi} \bar{\mathcal{J}})(i)|}{\xi(i)} \leq e \left( \max_{j=1 \dots n} \frac{|\mathcal{J}(j) - \bar{\mathcal{J}}(j)|}{\xi(j)} \right)$$

This inequality holds for any state  $i$ :

$$\max_{i=1 \dots n} \frac{|(T_{\pi} \mathcal{J})(i) - (T_{\pi} \bar{\mathcal{J}})(i)|}{\xi(i)} \leq e \left( \max_{j=1 \dots n} \frac{|\mathcal{J}(j) - \bar{\mathcal{J}}(j)|}{\xi(j)} \right)$$

$$\|T_{\pi} \mathcal{J} - T_{\pi} \bar{\mathcal{J}}\|_{\xi} \leq e \|\mathcal{J} - \bar{\mathcal{J}}\|_{\xi}$$

So,  $T_{\pi}$  is a contraction w.r.t.  $\|\cdot\|_{\xi}$

with modulus  $e$ .

$$\begin{aligned} \mathcal{J} \mathcal{J}(i) &= \min_c \sum_j P_{ij}(c) (\mathcal{J}(j) + \tau(i, j)) \\ (T_{\pi} \mathcal{J})(i) &= \sum_j P_{ij}(\pi(i)) (\mathcal{J}(j) + \tau(i, j)) \end{aligned}$$

Next: to show that  $T$  is a contraction wrt.  $\|\cdot\|_{\xi}$

from (I), We have shown already that

$$(T_{\pi} \mathcal{J})(i) \leq (T_{\pi} \bar{\mathcal{J}})(i) + \epsilon \xi(i) \max_{\delta=1..n} \frac{|\mathcal{J}(\delta) - \bar{\mathcal{J}}(\delta)|}{\xi(\delta)}$$

Take minimum over  $\pi$  on both sides to obtain

$$(T \mathcal{J})(i) \leq (T \bar{\mathcal{J}})(i) + \epsilon \xi(i) \max_{\delta=1..n} \frac{|\mathcal{J}(\delta) - \bar{\mathcal{J}}(\delta)|}{\xi(\delta)}$$

Interchange  $\mathcal{J}$  &  $\bar{\mathcal{J}}$  to obtain

$$(T \bar{\mathcal{J}})(i) \leq (T \mathcal{J})(i) + \epsilon \xi(i) \max_{\delta=1..n} \frac{|\mathcal{J}(\delta) - \bar{\mathcal{J}}(\delta)|}{\xi(\delta)}$$

$$|(T \mathcal{J})(i) - (T \bar{\mathcal{J}})(i)| \leq \epsilon \xi(i) \max_{\delta=1..n} \frac{|\mathcal{J}(\delta) - \bar{\mathcal{J}}(\delta)|}{\xi(\delta)}$$

$$\max_{i=1..n} \frac{|(T \mathcal{J})(i) - (T \bar{\mathcal{J}})(i)|}{\xi(i)} \leq \epsilon \max_{\delta=1..n} \frac{|\mathcal{J}(\delta) - \bar{\mathcal{J}}(\delta)|}{\xi(\delta)}$$

$$\text{Or, } \|T \mathcal{J} - T \bar{\mathcal{J}}\|_{\xi} \leq \epsilon \|\mathcal{J} - \bar{\mathcal{J}}\|_{\xi}$$

Thus,  $T$  is a contraction wrt  $\|\cdot\|_{\xi}$  with modulus  $\epsilon$ .



## Value Iteration (VI):

1) Choose some  $J_0$

2) Repeatedly apply  $T$

$$J_0 \xrightarrow{T} J_1 \xrightarrow{T} J_2 \xrightarrow{T} \dots$$

$$J_{k+1} = T J_k$$

For the "all policies are proper" case,

$$\|J_k - J^*\|_{\infty} = \|T^k J_0 - T^k J^*\|_{\infty} \leq e^k \|J_0 - J^*\|_{\infty}$$

$$\|T J - T \bar{J}\|_{\infty} \leq e \|J - \bar{J}\|_{\infty}$$

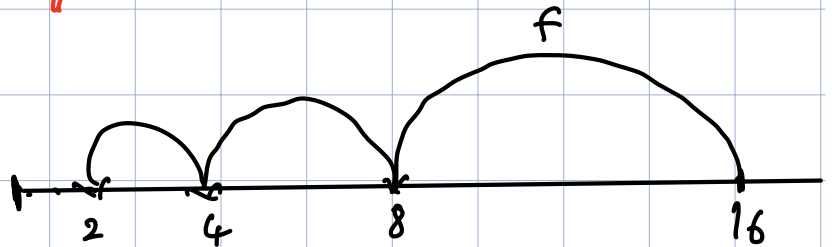
$$\Rightarrow \|T^k J - T^k \bar{J}\|_{\infty} \leq e^k \|J - \bar{J}\|_{\infty}$$

$\rightarrow$  VI error bound if all policies are proper

## Lecture-12\*

Contraction mappings: A quick detour

$$f(x) = \frac{x}{2}$$



For any  $x$ ,  $\lim_{n \rightarrow \infty} f^n(x) = 0$ ,

$$f(0) = 0$$

fixed point

Repeated application of  $f$  leads to the fixed point.

Vector space  $X$ .

Norm  $\|\cdot\|$  satisfies 3 properties:

(i)  $\|x\| = 0$  iff  $x = 0$

(ii)  $\|x + y\| \leq \|x\| + \|y\|$ ,  $\forall x, y \in X$ . "Triangle inequality"

(iii)  $\|cx\| = |c| \|x\|$ ,  $\forall x \in X$  & scalar  $c$ .

Contraction mapping:

$F: X \rightarrow X$  is a contraction mapping if

$\exists \rho \in (0, 1)$  s.t.

$\|F(x) - F(y)\| \leq \rho \|x - y\|$ ,  $\forall x, y \in X$ .

The space  $X$  is complete under the norm  $\|\cdot\|$  if every Cauchy sequence  $\{x_k\} \subset X$  converges.

A sequence  $\{x_k\}$  is Cauchy if  $\forall \epsilon > 0 \exists N_\epsilon$  s.t.

$\|x_m - x_n\| \leq \epsilon \quad \forall m, n \geq N_\epsilon$

$\forall \epsilon > 0, \exists N$  s.t.  $\|x_n - x^*\| < \epsilon \quad \forall n > N$  |  $\lim_{n \rightarrow \infty} x_n = x^*$



Fact (i) If  $X$  is complete &  $F$  is a contraction wrt  $\|\cdot\|$  with modulus  $\rho$ , then

" $F$  has a unique fixed point i.e.  $\exists$  an  $x^*$  s.t.  $F(x^*) = x^*$ "

(ii)  $x_{k+1} = F(x_k)$ . Then  $x_k \rightarrow x^*$  as  $k \rightarrow \infty$ .

$\{\frac{1}{n}\}$   $\lim_{n \rightarrow \infty} \frac{1}{n} = 0$

If the underlying space is restricted to  $(0, 1]$ , the  $\{\frac{1}{n}\}$  is Cauchy, but doesn't converge

# Contraction mappings in MDPs: (Ref: Sec 1.5 of DPOC-Vol II)

Weighted  
max-norm

$$\|T\|_{\xi} = \max_{x \in \mathcal{X}} \frac{|T(x)|}{\xi(x)}$$

where  $\xi(x) > 0 \quad \forall x \in \mathcal{X}$ .

If  $\xi(x) = 1 \quad \forall x$   
then,  $\|T\|_{\xi} = \|T\|_{\infty}$

$\|\cdot\|_{\infty} \rightarrow$  sup-norm or  
max-norm

Let  $B(\mathcal{X})$  denote all functions  $T$  s.t.  $\|T\|_{\xi} < \infty$ .

Simple case: finite state  
space.

Prop 3: Let  $F: B(\mathcal{X}) \rightarrow B(\mathcal{X})$  be a contraction

mapping with modulus  $\rho \in (0, 1)$ .

Assume  $B(\mathcal{X})$  is complete. Then,

(i) There exists a unique  $T^* \in B(\mathcal{X})$  s.t.

$$T^* = F T^*$$

$\forall I$  converges  
(asymptotically)

(ii)  $\lim_{k \rightarrow \infty} F^k T_0 = T^*$ , for any  $T_0 \in B(\mathcal{X})$

$$\text{Also, } \|F^k T_0 - T^*\|_{\xi} \leq \rho^k \|T_0 - T^*\|_{\xi}.$$

error bound for  $\forall I$   
that holds  $\forall k \geq 1$

Pf: See next page

Fix some  $T_0 \in B(X)$ .

Do  $T_{k+1} = FT_k$  starting with  $T_0$

$$\|T_{k+1} - T_k\|_{\mathcal{L}} = \|FT_k - FT_{k-1}\|_{\mathcal{L}} \leq e \|T_k - T_{k-1}\|_{\mathcal{L}}$$

↑  
F is a  $e$ -contraction

$$\Rightarrow \|T_{k+1} - T_k\|_{\mathcal{L}} \leq e^k \|T_1 - T_0\|_{\mathcal{L}} \quad (*)$$

So,  $\forall k \geq 0, m \geq 1$ , we have

$$\|T_{k+m} - T_k\|_{\mathcal{L}} = \left\| \underbrace{(T_{k+m} - T_{k+m-1})}_{\text{Telescoping sum}} + (T_{k+m-1} - T_{k+m-2}) + \dots + (T_{k+1} - T_k) \right\|_{\mathcal{L}}$$

$$\stackrel{\text{triangle inequality}}{\leq} \sum_{i=1}^m \|T_{k+i} - T_{k+i-1}\|_{\mathcal{L}}$$

$$\stackrel{\text{using } (*)}{\leq} e^k \underbrace{(1 + e + e^2 + \dots + e^{m-1})}_{\leq \sum_{i=0}^{\infty} e^i = \frac{1}{1-e}} \|T_1 - T_0\|_{\mathcal{L}}$$

$$\|T_{k+m} - T_k\|_{\mathcal{L}} \leq \frac{e^k}{1-e} \|T_1 - T_0\|_{\mathcal{L}} \quad (**)$$

Is  $\{T_k\}$  a Cauchy sequence? Yes. (note:  $\forall \epsilon \leq e$  in (\*\*))

Since  $B(X)$  is complete,  $T_k \rightarrow T^*$  and  $T^* \in B(X)$ .

↑  
 $T^*$  is only a limit of  $\forall \epsilon$  & we haven't shown  $FT^* = T^*$  yet.

To show:  $J^*$  is a fixed point.

$$\begin{aligned}\|FJ^* - J^*\|_{\xi} &\leq \|FJ^* - J_k\|_{\xi} + \|J_k - J^*\|_{\xi} \\ &\leq e \|J^* - J_{k-1}\|_{\xi} + \|J_k - J^*\|_{\xi} \\ &\xrightarrow{k \rightarrow \infty} 0 \quad \text{since } J_k \rightarrow J^*\end{aligned}$$

So,  $FJ^* = J^*$ .

To show: uniqueness of  $J^*$ .

Assume  $\tilde{J}$  is another fixed point ( $F\tilde{J} = \tilde{J}$ ).

$$\|J^* - \tilde{J}\|_{\xi} = \|FJ^* - F\tilde{J}\|_{\xi}$$

$$\|J^* - \tilde{J}\|_{\xi} \leq e \|J^* - \tilde{J}\|_{\xi}$$

$e \downarrow$   
 $e \in (0,1)$

$$\Rightarrow J^* = \tilde{J}$$

To show: error bound

$$\|F^k J_0 - J^*\|_{\xi} = \|F^k J_0 - FJ^*\|_{\xi}$$

$$\leq e \|F^{k-1} J_0 - J^*\|_{\xi}$$

& repeat to infer

$$\|F^k J_0 - J^*\|_{\xi} \leq e^k \|J_0 - J^*\|_{\xi}$$

Back to SSPs:

Assuming all policies are proper,  $\exists \beta$  s.t.  $\beta(i) > 0 \forall i$  and

$$\|TJ - T\bar{J}\|_{\infty} \leq \rho \|J - \bar{J}\|_{\infty}, \quad \forall J, \bar{J} \in \mathbb{R}^n$$

Using the error bound from the second claim in Prop 3, we obtain

$$\|J_k - J^*\|_{\infty} \leq \rho^k \|J_0 - J^*\|_{\infty},$$

where  $J_k = T J_{k-1}$  &  $J_0$  is the initial vector for value iteration.

Remark: A similar error bound holds for a evaluating a proper policy  $\pi$  using VI.

Lecture-13\*

Gauss-Seidel variant of VI:

Note: In VI, we do  $J_{k+1} = TJ_k$ , i.e., we apply the operator  $T$  for all states  $i$

Alternative: update one state at a time & we recent updates of states that are already updated.



## Gauss-Seidel VI update:

$$(FJ)(i) = \min_{a \in \mathcal{A}(i)} \sum_j P_{ij}(a) (g(i, a, j) + J(j))$$

same as  $(TJ)(i)$

For  $i = 2, \dots, n$ ,

$$(FJ)(i) = \min_{a \in \mathcal{A}(i)} \left[ \sum_j P_{ij}(a) g(i, a, j) + \sum_{j=1}^{i-1} P_{ij}(a) (FJ)(j) + \sum_{j=i}^n P_{ij}(a) J(j) \right]$$

For  $j=1 \dots i-1$  we previous updates

Remark: Gauss-Seidel VI converges faster than regular VI

$$\|J_k^{VI} - J^*\|_{\infty} \geq \|J_k^{GS-VI} - J^*\|_{\infty} \leftarrow \text{proof in Prop 2.2.5 of Vol. II.}$$

## Asynchronous VI

Synchronous VI:  $\forall i, \forall k, J_{k+1}(i) = TJ_k(i)$

(i) Start with an arbitrary  $J_0$

(ii) In  $k$ th iteration, pick an index  $i_k \in \{1, \dots, n\}$  and do

$$J_{k+1}(i) = \begin{cases} (TJ_k)(i) & \text{if } i = i_k \\ J_k(i) & \text{else} \end{cases}$$

If all states are picked infinitely often, then

$$J_k \rightarrow J^* \text{ as } k \rightarrow \infty$$

(Proof: Check Vol. II some section)

### Lecture 13\*

## Q-learning as a form of VI

When the model is known, i.e., we can compute  $T$ ,  
Q-learning is equivalent to VI.  
Cumulative Q-learning doesn't require the model. But we are following DROC Vol II nomenclature.

However, in the case when  $T$  is not computable directly, such as in a typical RL setting, there is a natural extension of "Q-learning" available.

Consider an SSP

$$VI: J_{k+1} = T J_k, \text{ with fixed } J_0$$

$$J_{k+1}(i) = \min_{a \in A(i)} \sum_j p_{ij}(a) (g(i, a, j) + J_k(j))$$

$$J_{k+1}(i) = \min_{a \in A(i)} Q_{k+1}(i, a), \quad \text{--- (1)}$$

$$\text{where } Q_{k+1}(i, a) = \sum_j p_{ij}(a) (g(i, a, j) + \min_{b \in A(j)} Q_k(j, b)) \quad \text{--- (2)}$$

with

$$J_0(i) = \min_{a \in A(i)} Q_0(i, a) \quad \text{--- (3)}$$

① - ③: Q-learning update

Why Q-learning works?

Q-factors / Q-values

Define  $Q^*(i, a) = \sum_j P_{ij}(a) (g(i, a, j) + J^*(j))$

Take action  $a$  in state  $i$  & from the next state onwards follow the optimal policy.

Bellman Equation ( $J^*(i) = T J^*(i), \forall i \Rightarrow J^*(i) = \min_a E(g(i, a, j) + J^*(j))$ )

$$J^*(i) = \min_a Q^*(i, a)$$

Thus,  $Q^*$  can be alternatively defined by

$$Q^*(i, a) = \sum_j P_{ij}(a) \left( g(i, a, j) + \min_{b \in A(j)} Q^*(j, b) \right)$$

Q-Bellman Equation ( $\Rightarrow Q^*(i, a) = E(g(i, a, j)) + \min_b Q^*(j, b)$ )

The operator underlying Q-Bellman equation is

$$(FQ)(i, a) = \sum_j P_{ij}(a) \left( g(i, a, j) + \min_{b \in A(j)} Q(j, b) \right)$$

It can be shown that  $FQ$  is a contraction mapping when all policies are proper

So, VI for Q-Bellman equation is start with  $Q_0(\cdot, \cdot)$  & keep applying  $(FQ)(\cdot, \cdot)$  operator

$$Q_0 \xrightarrow{F} Q_1 \xrightarrow{F} \dots \xrightarrow{\uparrow \text{eventually}} Q^*$$

Once you have  $Q^*$ , we can obtain  $J^*$  wrt  
 $J^*(i) = \min_a Q^*(i, a)$

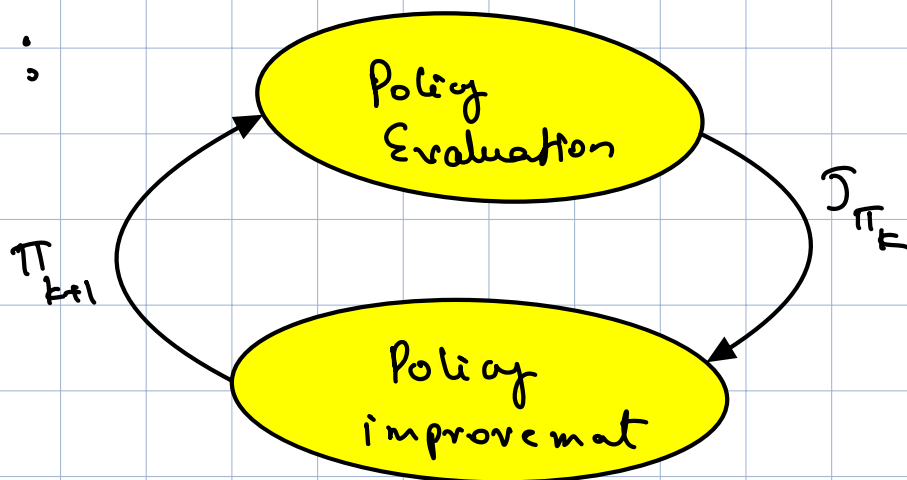
Eq (1)-(3) is just VI on Q-values.

# Policy iteration (PI)

VI can possibly take infinite # of iterations to converge

PI is a method "guaranteed" to converge within a finite # of iterations, for finite state-action spaced MDPs.

PI:



$$J_{\pi_0} \geq J_{\pi_1} \text{ -----}$$

↑  
ineq strict  
for at least  
one state

, if not, i.e.,  $J_{\pi_k} = J_{\pi_{k+1}}$ ,  
then, we have found the  
optimal policy.

"  $\exists i$  s.t.  $J_{\pi_0}(i) > J_{\pi_1}(i)$  "

PI algorithm:

Step 1: Start with a proper policy  $\pi_0$

Step 2: Evaluate  $\pi_k$ , i.e., compute  $J_{\pi}$   
(Policy Evaluation) by solving  $J = T_{\pi_k} J$

$$(\Rightarrow) J(i) = \sum_j P_{ij}(\pi_k(i)) (g(i, \pi_k(i), j) + J(j)), \quad \forall i$$

(here  $J(1) \dots J(n)$  are the unknowns & solving (\*) given  $J_{\pi_k}$ )

current policy  $\pi_k$

### Step 3: Policy improvement

Find a new policy  $\pi_{k+1}$  by

$$(\Rightarrow) \pi_{k+1}(i) = \arg \min_{a \in A(i)} \sum_j P_{ij}(a) (g(i, a, j) + J_{\pi_k}(j))$$

known (from step 2)

Bellman operator

next/improved policy  $\rightarrow \pi_{k+1}$

If  $J_{\pi_{k+1}}(i) < J_{\pi_k}(i)$  for at least one state  $i$ ,  
then go to step 2 & repeat.

Else, Stop. In this case,  $J_{\pi_{k+1}} = J_{\pi_k} = J^*$  (we will show this)

Remark: (Finding a proper policy to start with)

Step 1 can be modified by starting with a finite  $J$  & going to step 3 directly. That would give  $\pi_0$ , a proper policy.

This construction won't work for finding a proper policy

$(T_{\pi_0} J = J \Rightarrow \pi_0 \text{ is proper})$

And from there on do step 2 & 3 in tandem until convergence.

## Lecture-14

Claim: Policy improvement does just that.

Let  $\pi, \pi'$  be two proper policies s.t.

$$T_{\pi'} J_{\pi} = T J_{\pi}$$

Then,  $J_{\pi'}(i) \leq J_{\pi}(i) \quad \forall i \quad (*)$

with strict inequality for at least one of the states if  $\pi$  is not optimal.

PF: We know that  $J_{\pi} = T_{\pi} J_{\pi} \rightarrow J_{\pi}$  is a fixed point of  $T_{\pi}$

Also,  $T_{\pi'} J_{\pi} = T J_{\pi} \leftarrow$  given

So,

$$J_{\pi}(i) = \sum_j P_{ij}(\pi(i)) (g(i, \pi(i), j) + J_{\pi}(j))$$

$$\stackrel{T_{\pi'} J_{\pi} \geq T J_{\pi}}{\geq} \min_a \sum_j P_{ij}(a) (g(i, a, j) + J_{\pi}(j))$$

$$= (T_{\pi'} J_{\pi})(i), \quad \forall i$$

$$\text{eqn. (A)} \rightarrow J_{\pi} \geq T_{\pi'} J_{\pi} \geq T_{\pi'}^2 J_{\pi} \geq \dots \geq T_{\pi'}^k J_{\pi}$$

$$\Rightarrow J_{\pi} \geq \lim_{k \rightarrow \infty} T_{\pi'}^k J_{\pi} = J_{\pi'}$$

$\Rightarrow J_{\pi} \geq J_{\pi'} \quad$  So,  $(*)$  is proved.

To prove:  $\pi$  is optimal if the inequality in (\*\*) isn't strict for at least one state

Now, if  $J_{\pi} = J_{\pi'}$ , then from eqn (A),

$$J_{\pi} = T_{\pi'} J_{\pi} \quad \text{--- (1)}$$

$$T_{\pi'} J_{\pi} = T J_{\pi} \quad \text{by construction --- (2)}$$

Using (1) & (2), we get

$$J_{\pi} = T J_{\pi}$$

$\Rightarrow$   $\pi$  is optimal since  $T$  is the Bellman operator, which has a unique fixed point

Thus, if  $\pi$  isn't optimal, then  $J_{\pi'}(i) < J_{\pi}(i)$   
for at least one state  $i$



**Claim!** If the number of proper policies is finite, the PI converges in a finite number of steps.

Why? Policy improvement guarantees a better policy if the latter isn't optimal.



## Modified PI:

Let  $\{m_0, m_1, \dots\}$  be positive integers.

Let  $\mathcal{J}_1, \mathcal{J}_2, \dots$  and  $\pi_0, \pi_1, \dots$  be computed as follows:

Policy improvement  $\rightarrow T_{\pi_k} \mathcal{J}_k = T \mathcal{J}_k \rightarrow$  as in PI

Policy evaluation  $\rightarrow \mathcal{J}_{k+1} = T_{\pi_k}^{m_k} \mathcal{J}_k$   $\rightarrow$  Approx. policy evaluation by applying  $T_{\pi_k}$  " $m_k$ " times

by  $m_k$  steps of VI

Two special cases:

If  $m_k = \infty$ , Modified PI = regular PI  
since  $T_{\pi_k}^{\infty} \mathcal{J}_k = \mathcal{J}_{\pi_k}$

If  $m_k = 1$ , Modified PI = VI

$$\mathcal{J}_{k+1} = T_{\pi_k}^{m_k} \mathcal{J}_k = T_{\pi_k} \mathcal{J}_k = T \mathcal{J}_k$$

Recommended: Choose  $m_k > 1$

Convergence:

Modified PI converges

"proof skipped" See Ch. 2 of

Bertsekas DPOC-Vol II

## Asynchronous PI:

Let  $J_1, J_2, J_3, \dots$  : Sequence of optimal cost estimates  
 $\pi_1, \pi_2, \pi_3, \dots$  : Corresponding sequence of policies

Given  $(J_k, \pi_k)$ , we select a subset  $S_k$  of states and generate  $(J_{k+1}, \pi_{k+1})$  in one of the following two ways:

(Way 1)

$$J_{k+1}(i) = \begin{cases} (T_{\pi_k} J_k)(i) & \text{if } i \in S_k \\ J_k(i) & \text{else} \end{cases}$$

(+) *Possible iteration*

(Way 2)

$$\pi_{k+1}(i) = \begin{cases} \underset{a \in \mathcal{A}(i)}{\operatorname{argmin}} \left( \sum_j P_{ij}^a (g(i, a, j) + J_k(j)) \right), & \text{if } i \in S_k \\ \pi_k(i) & \text{else} \end{cases}$$

(+) *(+) is like policy improvement on subset  $S_k$*

(+)  *$\downarrow T J_k$*

## Special Cases of Async-PID:

- ① If  $S_k =$  entire state space,  $(*)$  is performed infinite # of times before one  $(+)$  operation, then we get **regular PID**
- ② If  $S_k =$  entire state space,  $(*)$  performed  $m_k$  times before  $(+)$ , then we get **modified PID**
- ③ If  $S_k =$  entire state space, one  $(*)$  operation followed immediately by one  $(+)$  operation, then we get **value iteration**
- ④ If  $|S_k| = 1$ , and one  $(*)$  operation followed immediately by one  $(+)$  operation, then we get **Asynchronous VI**

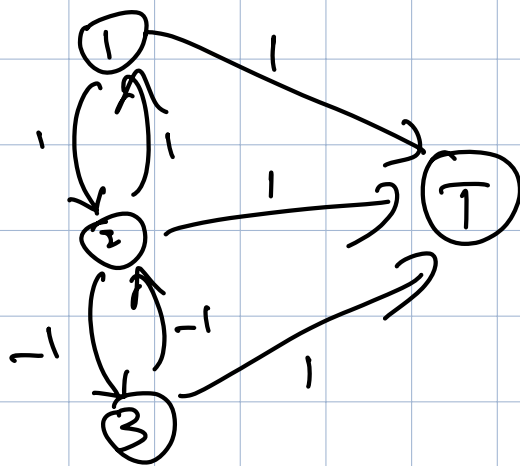
Remark! Async-PID can be shown to converge.  
See DPOC - Vol II for details.

Aside! Value iteration with Q-factors (leads finally to Q-learning (in an RL setting))  
PID leads to "Actor-critic methods".

## Discussion of some problems

① SSP  $\exists$  one proper policy  
Each improper  $\pi$  has  $J_{\pi}(i) = \infty$  for at least one  $i$

Claim:  $\nexists$  improper  $\pi'$  s.t.  $J_{\pi'}(i) = -\infty$  for some  $i$ .



$\rightarrow$  This is not a valid counterexample

Claim is indeed true.

Proof by contradiction.

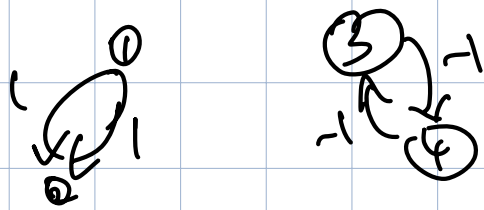
$\exists$  Proper  $\pi_1$ , &  $\exists$  improper  $\pi_2$  s.t.  
 $J_{\pi_2}(j) = -\infty$ .

Manufacture  $\pi_3$  that follows  $\pi_1$ ,

everywhere except  $j$ .

$\exists$  state  $i$  s.t.  $J_{\pi_3}(i) = \infty$ .  $\rightarrow$  violates assumption.

" $\exists$  one proper policy"  $\leftarrow$  why is this needed



(1)

(2)  $f$  is a contraction mapping with modulus  $\alpha$ .  
 $f(x^*) = x^*$

$$f(x) \leq x \Rightarrow x^* \leq x \quad \text{--- (1)}$$

One proof for (1):

$$f(f(x)) = f^2(x) \leq f(x) \leq x \quad \leftarrow \text{requires monotonicity}$$

$$\Rightarrow \lim_{k \rightarrow \infty} f^k(x) \leq x$$

$$\Rightarrow x^* \leq x$$

Monotonicity of  $f$ :  $x \leq y \Rightarrow f(x) \leq f(y)$

H.W. Think of a counterexample.

$$J_{N-1}(x) = \min_{\alpha} E( g( \quad ) + \underbrace{J_N(x')} )$$

$$\approx \min_{\alpha} E( g( \quad ) + J'_N(x') )$$

$$= J'_{N-1}(x)$$