

Lecture 14* (Contd)

Infinite horizon discounted MDPs

(Ref: Dpoc vol. II. Chapter 142)

Goal:

$$\mathcal{J}^*(i) = \min_{\pi \in \Pi} \mathcal{J}_\pi(i), \quad \forall i, \text{ where}$$

optimal discounted cost

$$\mathcal{J}_\pi(i) = E \left(\sum_{k=0}^{\infty} \alpha^k g(x_k, \pi(x_k), x_{k+1}) \mid x_0=i \right)$$

↓ discount
 $0 < \alpha < 1$
 ↗ single-stage cost
 ↗ policy
 ↗ start state

Let π^* denote the optimal policy i.e., $\arg \min_{\pi \in \Pi} \mathcal{J}_\pi(i)$

Let $S = \{1, \dots, n\}$ denote the state space

(A1) We shall assume the single stage cost is bounded, i.e.,

$$|g(\cdot, \cdot, \cdot)| \leq M < \infty.$$

Note: We do not assume existence of a special termination state.

Bellman and another operator:

For $J = (J(1), \dots, J(n))$,

define Bellman operator T as follows:

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

for a stationary policy π ,

$$(T_\pi \mathcal{J})(i) = \sum_{j=1}^n p_{ij}(\pi(i)) (g(i, \pi(i), j) + \alpha \mathcal{J}(j)), \quad \forall i$$

$$P_\pi = \begin{bmatrix} p_{11}(\pi(1)) \\ \vdots \\ p_{nn}(\pi(n)) \end{bmatrix}$$

$$g_\pi = \left(\begin{array}{c} \sum_{j=1}^n p_{1j}(\pi(1)) g(1, \pi(1), j) \\ \vdots \\ \sum_{j=1}^n p_{nj}(\pi(n)) g(n, \pi(n), j) \end{array} \right) = \begin{pmatrix} g(1, \pi(1)) \\ \vdots \\ g(n, \pi(n)) \end{pmatrix}$$

if single stage cost is
a function of
the current state &
current action, i.e.,
 $g(i, a, j) = g(i, a) \forall j$

$$\text{So, } T_\pi \mathcal{J} = g_\pi + \alpha P_\pi \mathcal{J}$$

Even in the discounted case, T, T_π are monotone.

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

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

(i) $(T^k \pi)(i) \leq (T^k \pi')(i)$, and

(ii) For any stationary policy π ,

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

Pf:

H.W.

The constant shift lemma holds here as well.

Lemma 2:

Stationary π , $\delta \rightarrow$ positive scalar, $e \rightarrow$ vector of n ones.

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

$$(i) (T^k(\pi + \delta e))(i) = (T^k \pi)(i) + \delta^k \delta$$

$$(ii) (T_\pi^k(\pi + \delta e))(i) = (T_\pi^k \pi)(i) + \delta^k \delta$$

Note: It is an equality here, while the corresponding SSP claim had an inequality (owing to terminal state T).

Pf:

H.W.

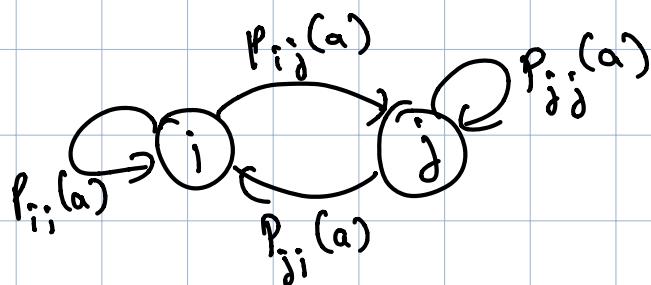
"Every discounted problem has an equivalent SSP".

Given discounted MDP on states $\{1, \dots, n\}$,

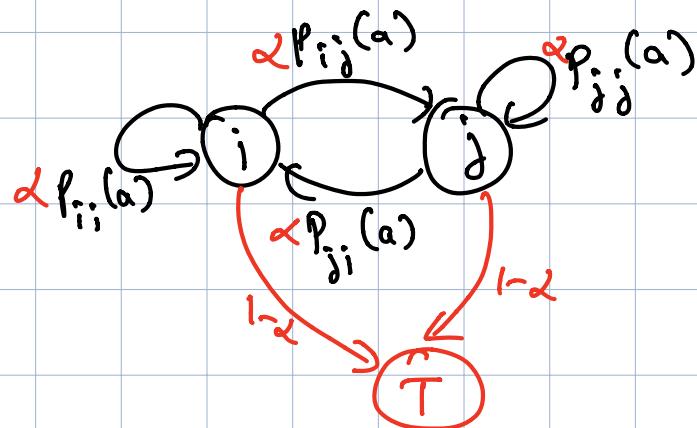
form an SSP on states $\{1, \dots, n\} \cup \{T\}$

\uparrow
add this extra
state & make it
lost-free & absorbing.

Idea: In the SSP, w.p. α pick a next state
according to transition probabilities of discounted MDP
& w.p. $(1-\alpha)$ move to "T".



discounted - MDP



Equivalent SSP

Why are these two MDPs equivalent?

① Note that in the SSP, all policies are proper.

② In the discounted MDP, the expected k-th stage cost
is $E(\alpha^k g(i, a, j)) = \alpha^k \sum_j p_{ij}(a) g(i, a, j)$

③ In the SSP, the expected k th stage cost is

$$\left[\sum_j p_{ij}(a) g(i, a, j) \right] \times \alpha^k$$

If the terminal state is not hit upto k th stage, then the underlying probabilities will have a α^k multiplier.

Lecture-15

Prop 1: (VI converges)

Assume (A1).

For any finite J , the optimal cost satisfies

$$J^*(i) = \lim_{N \rightarrow \infty} (T^N J)(i), \quad \forall i$$

(Corollary): For a stationary policy π , we have

$$J_\pi(i) = \lim_{N \rightarrow \infty} (T_\pi^N J)(i), \quad \forall i \text{ for any finite } J.$$

Pf) Given a policy $\pi = \{\mu_0, \mu_1, \dots\}$ and a state $i \in X$,

$$J_\pi(i) = \lim_{N \rightarrow \infty} E \left(\sum_{l=0}^{N-1} \alpha^l g(x_l, \mu_l(x_l), x_{l+1}) \right)$$

For convenience!
Drop the dependence
on i on RHS.

$$= E \left(\sum_{l=0}^{L-1} \alpha^l g(x_l, \mu_l(x_l), x_{l+1}) \right)$$

$$+ \lim_{N \rightarrow \infty} E \left(\sum_{l=L}^{N-1} \alpha^l g(x_l, \mu_l(x_l), x_{l+1}) \right)$$

term B

(1)

Since $|g(\cdot, \cdot, \cdot)| \leq M$ by assumption,

$$| \text{term B} | \leq M \sum_{l=L}^{\infty} \alpha^l = \frac{M \alpha^L}{1-\alpha} \quad (2)$$

$$E \left(\sum_{l=0}^{L-1} \alpha^l g(x_l, \mu_l(x_l), x_{l+1}) \right)$$

$$= \mathcal{T}_{\pi}(i) - \text{term B}$$

$$E \left(\alpha^L \mathcal{T}(x_L) + \sum_{l=0}^{L-1} \alpha^l g(x_l, \mu_l(x_l), x_{l+1}) \right)$$

$$= \mathcal{T}_{\pi}(i) - \text{term B} + E(\alpha^L \mathcal{T}(x_L))$$

Using (2),

$$\mathcal{T}_{\pi}(i) - \frac{M \alpha^L}{1-\alpha} - \alpha^L \max_{j \in X} |\mathcal{T}(j)|$$

$$\leq E \left(\alpha^L \mathcal{T}(x_L) + \sum_{l=0}^{L-1} \alpha^l g(x_l, \mu_l(x_l), x_{l+1}) \right)$$

$$\leq \mathcal{T}_{\pi}(i) + \frac{M \alpha^L}{1-\alpha} + \alpha^L \max_{j \in X} |\mathcal{T}(j)| - (3)$$

applied \mathcal{T}_{π} L times on \mathcal{T}

Taking minimum over π on all sides of ③, we obtain $\forall i \in X$ and any $L > 0$ that

$$\gamma^*(i) - \frac{M\alpha^L}{1-\alpha} - \alpha^L \max_{j \in X} |\gamma(j)|$$

$$\leq (\gamma^L \gamma)(i) \quad \xrightarrow{\text{because } \min_{\pi} \gamma_{\pi}^L = \gamma^L}$$

$$\leq \gamma^*(i) + \frac{M\alpha^L}{1-\alpha} + \alpha^L \max_{j \in X} |\gamma(j)| \quad - ④$$

Taking $L \rightarrow \infty$ on all sides of ④, we obtain
 (Note: $\alpha \in (0, 1)$)

$$\gamma^*(i) \leq \lim_{L \rightarrow \infty} (\gamma^L \gamma)(i) \leq \gamma^*(i)$$

& the claim follows. ■

Corollary: The claim follows by considering an MDP where the only feasible action in a state i is $\pi(i)$, $\forall i$, and invoking Prop 1.
 (Note $\gamma = \gamma_{\pi}$ for this MDP).

Prop 2: (Bellman equation)

The optimal discounted cost \mathcal{J}^* satisfies

$$\mathcal{J}^* = T \mathcal{J}^*, \text{ i.e.,}$$

$$\mathcal{J}^*(i) = \min_a \sum_j p_{ij}(a) (g(i, a, j) + \gamma \mathcal{J}^*(j))$$

Also, \mathcal{J}^* is the unique fixed point of T .

PF) Eq ④ in the proof above is

$$\mathcal{J}^*(i) - \frac{M\gamma^L}{1-\gamma} - \gamma^L \max_{j \in X} (\mathcal{J}(j))$$

$$\leq (T^L \mathcal{J})(i)$$

$$\leq \mathcal{J}^*(i) + \frac{M\gamma^L}{1-\gamma} + \gamma^L \max_{j \in X} |\mathcal{J}(j)| \quad \text{--- (5)}$$

Applying operator T on all sides,

$$T \mathcal{J}^*(i) - \frac{M\gamma^{L+1}}{1-\gamma} - \gamma^{L+1} \max_{j \in X} (\mathcal{J}(j)) \leq (T^{L+1} \mathcal{J})(i)$$

$$\leq T \mathcal{J}^*(i) + \frac{M\gamma^{L+1}}{1-\gamma} + \gamma^{L+1} \max_{j \in X} |\mathcal{J}(j)|$$

*we add
constant shift
term.*

Taking $L \rightarrow \infty$ on all sides of the equation above,

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

(Uniqueness): Let J' be another fixed point of T .

$$J' = T J' = T^2 J' = \dots = \lim_{L \rightarrow \infty} T^L J' = J^*$$

■

(Corollary): For a stationary policy π , the associated cost J_π satisfies

$$J_\pi = T_\pi J_\pi \quad (\text{or})$$

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

$\forall i$

Also, J_π is the unique fixed point of T_π .

(Pf) Follows from Prop 2.

■

Lecture-16*

Necessary & sufficient condition for optimal policy:

Prop 3: A stationary policy π is optimal if and only if

$\pi(i)$ attains the minimum in the Bellman equation,

$\forall i \in \mathcal{X}$. Or, equivalently,

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

Pf>

$$\text{Assume } T\mathcal{J}^* = T_\pi \mathcal{J}^*$$

$$\text{We know } \mathcal{J}^* = T\mathcal{J}^*$$

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

$$\Rightarrow \mathcal{J}^* = \mathcal{J}_\pi \Rightarrow \pi \text{ is optimal.}$$

Converse: π is optimal

$$\Rightarrow \mathcal{J}^* = \mathcal{J}_\pi \Rightarrow \mathcal{J}^* = T_\pi \mathcal{J}^*$$

$$\text{Also, Bellman equation} \Rightarrow \mathcal{J}^* = T\mathcal{J}^*$$

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



Contraction property of T and T_π :

$$\text{Max-norm: } \|\mathcal{J}\|_\infty = \max_{i=1 \dots n} |\mathcal{J}(i)|$$

We will show that T, T_π are α -contractions in $\|\cdot\|_\infty$.

Prop 4: For any two bounded π, π' , and $\forall k \geq 1$

$$\|T^k \pi - T^k \pi'\|_{\infty} \leq \alpha^k \|\pi - \pi'\|_{\infty} \quad (*)$$

Pf: Let $c = \max_{i=1 \dots n} |\pi(i) - \pi'(i)|$

$$\pi(i) - c \leq \pi'(i) \leq \pi(i) + c \quad ①$$

Apply T^k on all sides of ① & use
constant shift lemma to obtain

$$(T^k \pi)(i) - \alpha^k c \leq (T^k \pi')(i) \leq (T^k \pi)(i) + \alpha^k c$$

$$\Rightarrow |(T^k \pi)(i) - (T^k \pi')(i)| \leq \alpha^k c \quad \text{holds for any } i$$

Taking max over i implies $(*)$.

■

Corollary: For any stationary π & bounded π, π' , and $\forall k \geq 1$

$$\|T_{\pi}^k \pi - T_{\pi}^k \pi'\|_{\infty} \leq \alpha^k \|\pi - \pi'\|_{\infty}$$

Value Iteration! Start with π_0 & repeatedly apply T .

Error bound for VI:

$$\|T^k \pi_0 - \pi^*\|_{\infty} \leq \alpha^k \|\pi_0 - \pi^*\|_{\infty}$$

Why? Set $\pi^* = \pi^{**}$ in (*) & note $T^k \pi^{**} = \pi^{**}$.

Example: Machine replacement

Recall n -states $1 \dots n$

Operating Cost $g(i)$

$$g(1) \leq g(2) \leq \dots \leq g(n)$$

p_{ij} \rightarrow transition probabilities

actions: do nothing & repair ($\text{Repair cost } R$)

Goal: minimize infinite horizon discounted cost
($\alpha \leftarrow$ discount factor)

Bellman equation:

$$\pi^*(i) = \min \left\{ \underbrace{R + g(i) + \alpha \pi^*(i)}_{\text{repair}}, \right.$$

$$\left. g(i) + \alpha \sum_{j=1}^n p_{ij} \pi^*(j) \right\}$$

So, optimal action is to repair if

$$R + g(i) + \alpha \pi^*(i) < g(i) + \alpha \sum_{j=1}^n p_{ij} \pi^*(j)$$

& "do nothing" otherwise.

Assume . ① $P_{ij} = 0$ if $j < i$ ← machine can't get better if you do nothing

$$\textcircled{2} \quad P_{ij} \leq P_{(i+1)j} \text{ if } i < j$$

Suppose \mathcal{T} is monotone non-decreasing, i.e.,
 $\mathcal{T}(1) \leq \mathcal{T}(2) - \dots \leq \mathcal{T}(n)$

Then,

$$\sum_{j=1}^n P_{ij} \mathcal{T}(j) \leq \sum_{j=i}^n P_{(i+1)j} \mathcal{T}(j), \quad i=1 \dots n-1$$

Since $g(i)$ is non-decreasing, we have

$(T\mathcal{T})(i)$ is non-decreasing in i , if \mathcal{T} is
 non-decreasing.

$\Rightarrow (T^k \mathcal{T})(i)$ is non-decreasing in i , $\forall k$

So, $\mathcal{T}^*(i) = \lim_{k \rightarrow \infty} (T^k \mathcal{T})(i)$ is non-decreasing in i

So, the function $\left(g(i) + 2 \sum_j P_{ij} \mathcal{T}^*(j) \right)$
 is non-decreasing in i .

Set of states $S_R = \left\{ i \mid R + g(i) + 2 \sum_j P_{ij} \mathcal{T}^*(j) \leq g(i) + 2 \sum_j P_{ij} \mathcal{T}^*(j) \right\}$

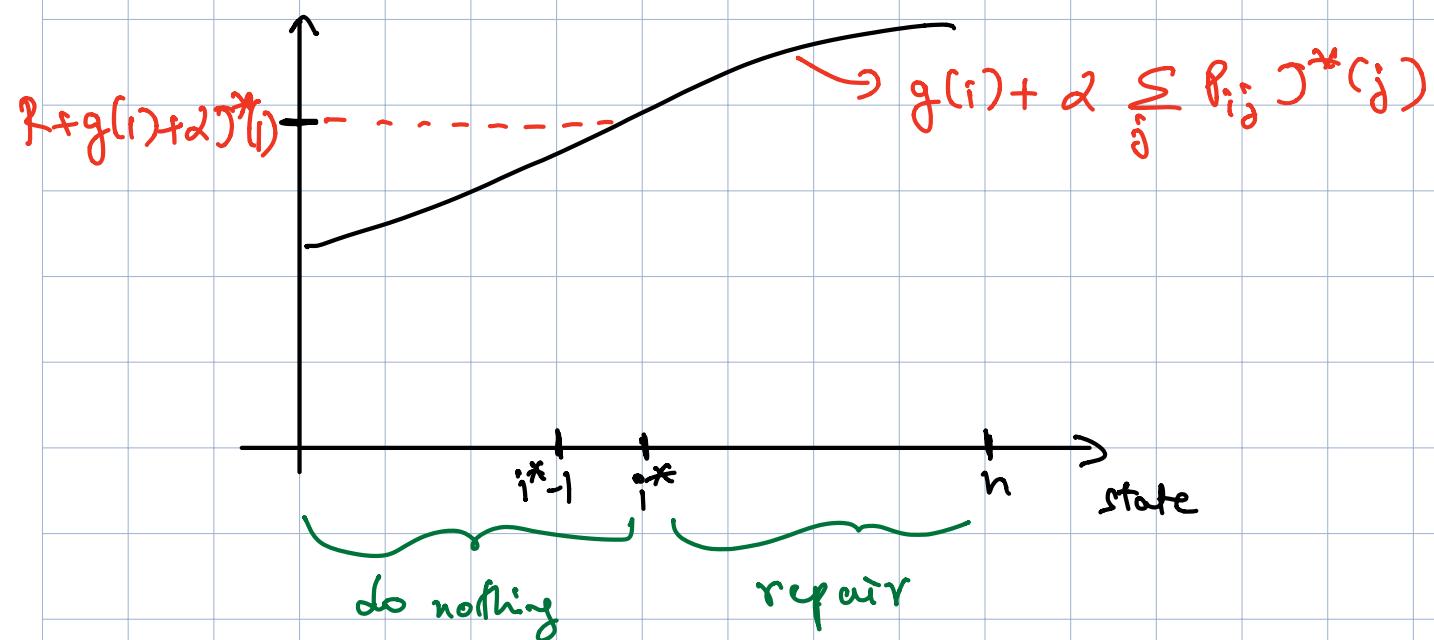
S_R = set of states where it is optimal to repair

$$i^* = \begin{cases} \text{smallest state in } S_R & \text{if } S_R \neq \emptyset \\ n+1 & \text{else} \end{cases}$$

Optimal policy

$$= \begin{cases} \text{repair} & \text{if } i \geq i^* \\ \text{do nothing} & \text{else} \end{cases}$$

A threshold-based
optimal policy



H.W. Think about policy iteration for this problem.

In particular, if we start with a threshold-based policy & do policy improvement, then does it lead to another threshold policy?

If yes, then π^* converges to optimal policy in at most n iterations.

Illustrative example for VI:

$$\text{MDP} \quad S = \{1, 2\} \quad A = \{a, b\}$$

$$P(a) = \begin{bmatrix} P_{11}(a) & P_{12}(a) \\ P_{21}(a) & P_{22}(a) \end{bmatrix} = \begin{bmatrix} \frac{3}{4} & \frac{1}{4} \\ \frac{3}{4} & \frac{1}{4} \end{bmatrix}$$

$$P(b) = \begin{bmatrix} P_{11}(b) & P_{12}(b) \\ P_{21}(b) & P_{22}(b) \end{bmatrix} = \begin{bmatrix} \frac{1}{4} & \frac{3}{4} \\ \frac{1}{4} & \frac{3}{4} \end{bmatrix}$$

Costs: $g(1, a) = 2, \quad g(1, b) = 0.5$
 $g(2, a) = 1, \quad g(2, b) = 3$

Discount $\alpha = 0.9$

$$\mathcal{J}_0 = (0, 0)$$

$$(\mathcal{T}\mathcal{J})(i) = \min \left\{ g(i, a) + \alpha \sum_{j=1}^2 P_{ij}(a) \mathcal{J}(j), \right.$$

$$\left. g(i, b) + \alpha \sum_{j=1}^2 P_{ij}(b) \mathcal{J}(j) \right\}$$

$$\mathcal{J}_1 = \mathcal{T}\mathcal{J}_0 = (0.5, 1)$$

$$\mathcal{J}_2 = (1.28, 1.56)$$

& so on.

PI algorithm:

Step 1: Start with a policy π_0

Step 2: Evaluate π_k , i.e., compute \mathcal{T}_π
(Policy Evaluation) by solving $\mathcal{T} = \mathcal{T}_{\pi_k} \mathcal{T}$

$$\Rightarrow \mathcal{T}(i) = \sum_j p_{ij}(\pi_k(i)) (g(i, \pi_k(i), j) + \gamma \mathcal{T}(j)), \quad \forall i$$

(here $\mathcal{T}(1), \dots, \mathcal{T}(n)$ are the unknowns of
solving (*) given \mathcal{T}_{π_k})

Step 3: Policy improvement

Find a new policy π_{k+1} by

$$\mathcal{T}_{\pi_{k+1}} \mathcal{T}_{\pi_k} = \mathcal{T} \mathcal{T}_{\pi_k}$$

$$\Rightarrow \pi_{k+1}(i) = \arg \min_{a \in A(i)} \sum_j p_{ij}(a) (g(i, a, j) + \gamma \mathcal{T}_{\pi_k}(j))$$

If $\mathcal{T}_{\pi_{k+1}}(i) < \mathcal{T}_{\pi_k}(i)$ for at least one state i ,

then go to step 2 & repeat.

Remark: Policy improvement claim holds even in
the discounted setting.

Policy improvement claim:

Let π, π' be two proper policies s.t.

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

Then,

$$T_{\pi'}(i) \leq T_{\pi}(i) \quad \forall i$$

with strict inequality for at least one of the states if π is not optimal.

Pf: Follows by a parallel argument to the proof in SSP case.

Lecture-17

PI example:

$$S = \{1, 2\}, A = \{a, b\}$$

$$P(a) = \begin{bmatrix} P_{11}(a) & P_{12}(a) \\ P_{21}(a) & P_{22}(a) \end{bmatrix} = \begin{bmatrix} 3/4 & 1/4 \\ 3/4 & 1/4 \end{bmatrix}$$

$$P(b) = \begin{bmatrix} P_{11}(b) & P_{12}(b) \\ P_{21}(b) & P_{22}(b) \end{bmatrix} = \begin{bmatrix} 1/4 & 3/4 \\ 1/4 & 3/4 \end{bmatrix}$$

$$\text{Costs: } g(1, a) = 2, \quad g(1, b) = 0.5$$

$$g(2, a) = 1, \quad g(2, b) = 3$$

$$\text{Discount } d = 0.9$$

Initialization!

$$\pi_0(1) = a, \quad \pi_0(2) = b$$

Policy evaluation:

Finding \mathcal{J}_{π_0} :

$$\mathcal{J}_{\pi_0}(1) = g(1, a) + \gamma P_{11}(a) \mathcal{J}_{\pi_0}(1) + \gamma P_{12}(a) \mathcal{J}_{\pi_0}(2)$$

$$\mathcal{J}_{\pi_0}(2) = g(2, b) + \gamma P_{21}(b) \mathcal{J}_{\pi_0}(1) + \gamma P_{22}(b) \mathcal{J}_{\pi_0}(2)$$

Using MDP data, we have

$$\mathcal{J}_{\pi_0}(1) = 2 + 0.9 \times \frac{3}{4} \times \mathcal{J}_{\pi_0}(1) + 0.9 \times \frac{1}{4} \times \mathcal{J}_{\pi_0}(2)$$

$$\mathcal{J}_{\pi_0}(2) = 3 + 0.9 \times \frac{1}{4} \times \mathcal{J}_{\pi_0}(1) + 0.9 \times \frac{3}{4} \times \mathcal{J}_{\pi_0}(2)$$

Solving,

$$\mathcal{J}_{\pi_0}(1) = 24.12, \quad \mathcal{J}_{\pi_0}(2) = 25.96$$

Policy improvement:

$$T_{\pi_1} \mathcal{J}_{\pi_0} = T \mathcal{J}_{\pi_0}$$

$$(T \mathcal{J}_{\pi_0})(1) = \min \left\{ 2 + 0.9 \left(\frac{3}{4} \times 24.12 + \frac{1}{4} \times 25.96 \right), \underbrace{\dots}_{\text{action } a} \right.$$

$$\left. 0.5 + 0.9 \left(\frac{1}{4} \times 24.12 + \frac{3}{4} \times 25.96 \right) \right\} \underbrace{\dots}_{\text{action } b}$$

$$= \min \{ 24.12, 23.45 \} = 23.45 \xrightarrow{\text{for action } b}$$

$$(TJ_{\pi_0})(2) = \min \left\{ 1 + 0.9 \left(\frac{3}{4} \times 24.12 + \frac{1}{4} \times 25.96 \right), \right.$$

action a

$$\left. 3 + 0.9 \left(\frac{1}{4} \times 24.12 + \frac{3}{4} \times 25.96 \right) \right\}$$

$$= \min \{ 23.12, 25.95 \} = 23.12 \quad \xleftarrow{\text{for action a}}$$

$$\pi_1(1) = b, \quad \pi_1(2) = a$$

Policy evaluation: J_{π_1} ?

$$J_{\pi_1}(1) = 7.33, \quad J_{\pi_1}(2) = 7.67$$

Policy improvement: $T_{\pi_2} J_{\pi_1} = T J_{\pi_1}$

$$\pi_2(1) = b, \quad \pi_2(2) = a$$

So, stop & output π_2

Linear programming

Want to solve $J^* = T J^*$

Idea: is to form a linear optimization problem whose solution is J^*

How? We know $\lim_{N \rightarrow \infty} T^N J = J^*$ for any J .

Suppose

Then,

$$J \leq T J$$

$$J \leq T^2 J$$

:

$$J \leq T^K J$$

\Rightarrow

$$J \leq J^* = T J^*$$



J^* is the largest solution that satisfies $J \leq T J$

Optimization problem:

constraint: $J \leq T J$

objective: $\max J$

More precisely, $J \leq T J$



$$J(i) \leq g(i, a) + \alpha \sum_{j=1}^n p_{ij}(a) J(j)$$

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

So, the LP formulation is:

Variables: $\lambda_1, \dots, \lambda_n$

Objective: $\max \sum_i \lambda_i$

Subject to

$$\lambda_i \leq g(i, a) + \alpha \sum_j p_{ij}(a) \lambda_j, \quad \text{for } i=1, \dots, n, a \in A(i)$$

(x)

Remark! Assuming $A(i) = \mathbb{A}$ & $|A| = q$,
we have $n \times q$ constraints in the LP (x)

variables = $n \leftarrow$ Cardinality of the state space

On problems with a large state space, LP is not practical.

Remark! Can we LP approach for solving SSPs as well.

H.W.: Write down the LP for the 2-state 2-action example used for VI/PI above.