

Infinite horizon discounted MDPs

(Ref: DPOC vol. II, Chapter 142)

Goal:

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

optimal discounted cost

Expected cumulative discounted cost aka value function

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

Single-stage cost

discount

policy

start state

 $0 < \alpha < 1 \rightarrow$ useful in financial applications as future costs/rewards are discounted.Let π^* denote the optimal policy i.e., $\arg \min_{\pi \in \Pi} J_{\pi}(i)$ (A1) The single stage cost $|g(x, a, x')| \leq M < \infty \quad \forall x, x' \in \mathcal{X}, a \in \mathcal{A}$

$$(A1) + (0 < \alpha < 1) \Rightarrow |J_{\pi}(i)| \leq M \sum_{k=0}^{\infty} \alpha^k = \frac{M}{1-\alpha} < \infty$$

"Unlike SSPs, we do not require the existence of a terminal state".

Bellman and another operatorFor $J = (J(1), \dots, J(n))$,define Bellman operator T as follows:

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

for a stationary policy π ,

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

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

array
 $X = \{1, \dots, n\}$

$$g_{\pi} = \begin{pmatrix} \sum_{j=1}^n P_{1j}(\pi(1)) g(i, \pi(1), j) \\ \vdots \\ \sum_{j=1}^n P_{nj}(\pi(n)) g(n, \pi(n), j) \end{pmatrix}$$

with this notation,

$$T_{\pi} J = g_{\pi} + \alpha P_{\pi} J$$

Remark: If single stage cost is a function of current state & action, i.e., $g(i, a)$, then

$$g_{\pi} = \begin{pmatrix} g(1, \pi(1)) \\ \vdots \\ g(n, \pi(n)) \end{pmatrix}$$

Even in the discounted case, T, T_π 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) $(T^k J)(i) \leq (T^k J')(i)$, and

(ii) For any stationary policy π ,

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

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

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

Pf: H.W.

"Every discounted problem has an equivalent SSP".

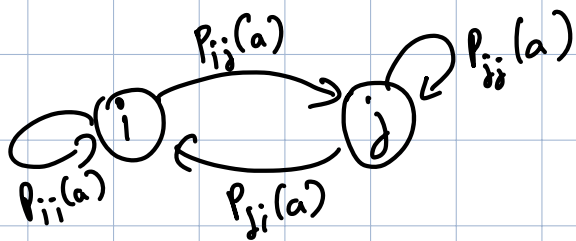
Given discounted MDP on states $\{1, \dots, n\}$,
form an SSP on states $\{1, \dots, n\} \cup \{T\}$



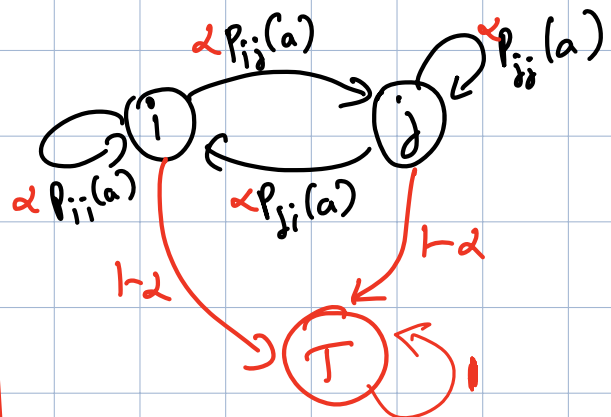
add this extra state & make it cost-free & absorbing.

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" & incurs no cost



Original discounted MDP



"Equivalent SSP"

$$g(i, a, j) = \begin{cases} \frac{g(i, a, j)}{\alpha} & \text{if } j \neq T \\ 0 & \text{else} \end{cases}$$

This scaling isn't necessary if cost g is of the form $g(i, a)$, i.e., indep. of next state

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.

"Optimal value is the same for the discounted MDP & the equivalent SSP"

Let \vec{J}^* → optimal value in SSP, Let $\vec{P}_{ij}(a)$ denote transition prob. here
 \vec{J}^* → optimal value in discounted MDP

Bellman equation in SSP

$$\vec{J}^*(i) = \min_a \sum_{j \in S} \vec{P}_{ij}(a) (g(i, a, j) + \vec{J}^*(j))$$

$$= \min_a \sum_{j \in S} \vec{P}_{ij}(a) g(i, a, j) + \sum_{j \in S} \vec{P}_{ij}(a) \vec{J}^*(j) + \vec{P}_{iT}(a) \vec{J}^*(T)$$

→ $S = \{1, \dots, n\}$

$$= \min_a \sum_{j \in S} \vec{P}_{ij}(a) g(i, a, j) + \sum_{j \in S} \alpha \vec{P}_{ij}(a) \vec{J}^*(j) + (1-\alpha) \vec{J}^*(T)$$

→ 0

$$\vec{J}^*(i) = \min_a \sum_{j \in S} \alpha \vec{P}_{ij}(a) \frac{g(i, a, j)}{\alpha} + (1-\alpha) g(i, a, T) + \alpha \sum_{j \in S} \vec{P}_{ij}(a) \vec{J}^*(j) \quad \text{--- } (\infty)$$

→ 0

Bellman equation in discounted MDP

$J^* = T J^*$ for discounted MDP

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

So, from (*) $J^*(i) = T J^*(i)$ (among $J = T J$ has a unique fixed point \rightarrow will be shown next)
 \Rightarrow optimal values coincide.

Lecture-16*

Prop 1: (VI converges)

Assume (A1).

For any finite J , the optimal cost satisfies

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

(Corollary: For a stationary policy π , we have

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

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

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

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

\rightarrow For convenience, drop the dependence on i on RHS.

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

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

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

$$-\frac{M \alpha^L}{1-\alpha} \leq \text{term B} \leq \frac{M \alpha^L}{1-\alpha}$$

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

$$= J_{\pi}(i) - \text{term B}$$

$$\begin{aligned} A &= J_{\pi} + C \\ |C| &\leq M \\ J_{\pi} - M &\leq A \leq J_{\pi} + M \end{aligned}$$

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

$$= J_{\pi}(i) - \text{term B} + E(\alpha^L J(x_L))$$

Using (2),

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

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

$$\leq J_{\pi}(i) + \frac{M \alpha^L}{1-\alpha} + \alpha^L \max_{j \in X} |J(j)| \quad (3)$$

applying J_{π} L times on J i.e., $\rightarrow = J_{\pi}^L J$

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

$$\begin{aligned}
 J^*(i) &= \frac{M\alpha^L}{1-\alpha} - \alpha^L \max_{j \in X} |J(j)| \\
 &\leq (T^L J)(i) \quad \longrightarrow \text{because } \min_{\pi} T_{\pi}^L = T^L \\
 &\leq J^*(i) + \frac{M\alpha^L}{1-\alpha} + \alpha^L \max_{j \in X} |J(j)| \quad - (4)
 \end{aligned}$$

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

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

& 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 $T = T_{\pi}$ for this MDP).

Prop 2: (Bellman equation)

The optimal discounted cost J^* satisfies

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

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

Bellman equation \rightarrow

Also, J^* is the unique fixed point of T .

PF) Eq (4) in the proof above is

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

$$\leq (T^L J)(i)$$

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

Applying operator T on all sides,

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

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

\rightarrow we used constant shift lemma.

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

$$T J^*(i) = J^*(i), \quad \forall i. \Rightarrow J^* \text{ is a fixed point of } T$$

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

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

$\therefore T$ has a unique fixed point \blacksquare

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

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

$$J_\pi(i) = \sum_j P_{ij}(\pi(i)) \left(g(i, \pi(i), j) + \alpha J_\pi(j) \right), \quad \forall i$$

Also, J_π is the unique fixed point of T_π .

Pf) Follows from Prop 2. \blacksquare

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,

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

Pf \Rightarrow Assume $T J^* = T_{\pi} J^*$ — (1)
We know $J^* = T J^*$ — (2)
 $J^* = T_{\pi} J^*$
 $\Rightarrow J^* = J_{\pi} \Rightarrow \pi$ is optimal

Converse: π is optimal

$$\Rightarrow J^* = J_{\pi} \Rightarrow J^* = T_{\pi} J^* \text{ — (1')}$$

$$\text{From BE, } J^* = T J^* \text{ — (2')}$$

$$\text{(1)' + (2)' } \Rightarrow T_{\pi} J^* = T J^* \quad \blacksquare$$

Contraction property of T and T_{π} :

$$\text{Max-norm: } \|J\|_{\infty} = \max_{i \in \mathcal{X}} |J(i)|$$

We will show that T, T_{π} are α -contractions in $\|\cdot\|_{\infty}$.

$\Rightarrow T$ is a contraction in $\|\cdot\|_\infty$ -norm with modulus α

Prop 4: For any two bounded $\mathcal{J}, \mathcal{J}'$, and $\forall k \geq 1$

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

\hookrightarrow modulus of contraction

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

$$\mathcal{J}(i) - c \leq \mathcal{J}'(i) \leq \mathcal{J}(i) + c \quad \text{--- (1)}$$

Apply T "k" times on all sides of (1) to get

holds $\forall i \rightarrow (T^k \mathcal{J})(i) - \alpha^k c \leq (T^k \mathcal{J}')(i) \leq (T^k \mathcal{J})(i) + \alpha^k c$

$$\Rightarrow |(T^k \mathcal{J})(i) - (T^k \mathcal{J}')(i)| \leq \alpha^k c, \forall i$$

$$\max_{i=1 \dots n} |(T^k \mathcal{J})(i) - (T^k \mathcal{J}')(i)| \leq \alpha^k c$$

$$\Rightarrow \|T^k \mathcal{J} - T^k \mathcal{J}'\|_\infty \leq \alpha^k \|\mathcal{J} - \mathcal{J}'\|_\infty$$

■

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

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

Value Iteration: Start with \mathcal{J}_0 & repeatedly apply T .

Error-bound for VI:

$$\|T^k \mathcal{J}_0 - \mathcal{J}^*\|_\infty \leq \alpha^k \|\mathcal{J}_0 - \mathcal{J}^*\|_\infty$$

Why? Set $J^1 = J^*$ in (*) & note $T^k J^* = J^*$.

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 (do nothing action)

actions: do nothing & repair (Repair cost R)

Goal: minimize infinite horizon discounted cost
(α discount factor)

Bellman equation: $J^*(i) = T J^*(i) = \min_a E(g(i,a,j) + \alpha J^*(j))$

$$J^*(i) = \min \left\{ \underbrace{R + g(i) + \alpha J^*(i)}_{\text{repair}}, \underbrace{g(i) + \alpha \sum_{j=1}^n P_{ij} J^*(j)}_{\text{do nothing}} \right\}$$

Optimal action: repair if

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

& "do nothing" otherwise

Assume (B1) $P_{ij} = 0$ if $j < i$ ← machine won't get better if we don't repair

(B2) $P_{ij} \leq P_{(i+1)j}$ if $i < j$ ← e.g. $j=10$
 $P_{9,10} \leq P_{8,10}$

Suppose J is monotone non-decreasing, i.e.,
 $J(1) \leq J(2) \leq \dots \leq J(n)$

Then,

Because of (B2) & monotone J →
$$\sum_{j=1}^n P_{ij} J(j) \leq \sum_{j=1}^n P_{(i+1)j} J(j), \quad i=1, \dots, n-1$$

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

$(TJ)(i)$ is non-decreasing in i , if J is non-decreasing. since $TJ(i) = \min(R + g(i) + \alpha J^*(i), g(i) + \alpha \sum P_{ij} J^*(j))$

⇒ $(T^k J)(i)$ is non-decreasing in i , $\forall k$

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

So, the function $g(i) + \alpha \sum_{j \geq i} P_{ij} J^*(j)$ is non-decreasing in i

Set of states $S_R = \left\{ i \mid R + g(i) + \alpha J^*(i) \leq g(i) + \alpha \sum_{j \geq i} P_{ij} J^*(j) \right\}$

$$\sum_{\tilde{j}=1}^n P_{i\tilde{j}} \mathcal{T}(\tilde{j}) \leq \sum_{\tilde{j}=1}^n P_{(i+1)\tilde{j}} \mathcal{T}(\tilde{j}),$$

$$n=3, \quad i=1, \quad i+1=2 \quad \mathcal{T}(\tilde{j}) \geq 1.$$

$$P_{11} + P_{12} + P_{13} = 1 \quad P_{ij} = 0 \text{ if } \tilde{j} < i$$

$$\leq \cancel{P_{21}} + P_{22} + P_{23} = 1$$

$$P_{11} \mathcal{T}_1 + P_{12} \mathcal{T}_2 + P_{13} \mathcal{T}_3$$

$$\begin{aligned} \mathcal{T}_1 &\leq \mathcal{T}_2 \leq \mathcal{T}_3 \\ P_{11} + P_{12} + P_{13} &= 1 \\ P_{22} + P_{23} &= 1 \end{aligned}$$

$$\leq P_{22} \mathcal{T}_2 + P_{23} \mathcal{T}_3$$

$$P_{11} \mathcal{T}_1 + P_{12} \mathcal{T}_2 + P_{13} \mathcal{T}_3$$

$$\leq P_{11} \mathcal{T}_2 + P_{12} \mathcal{T}_2 + P_{13} \mathcal{T}_3$$

$$P_{11} \mathcal{T}_1 \leq (P_{22} - P_{12}) \mathcal{T}_2 + (P_{23} - P_{13}) \mathcal{T}_3$$

$$(P_{22} - P_{12}) \mathcal{T}_1 + (P_{23} - P_{13}) \mathcal{T}_1 \leq (P_{22} - P_{12}) \mathcal{T}_2 + (P_{23} - P_{13}) \mathcal{T}_3$$

$$P_{11} \mathcal{T}_1 \leq \text{RHS} \leftarrow$$

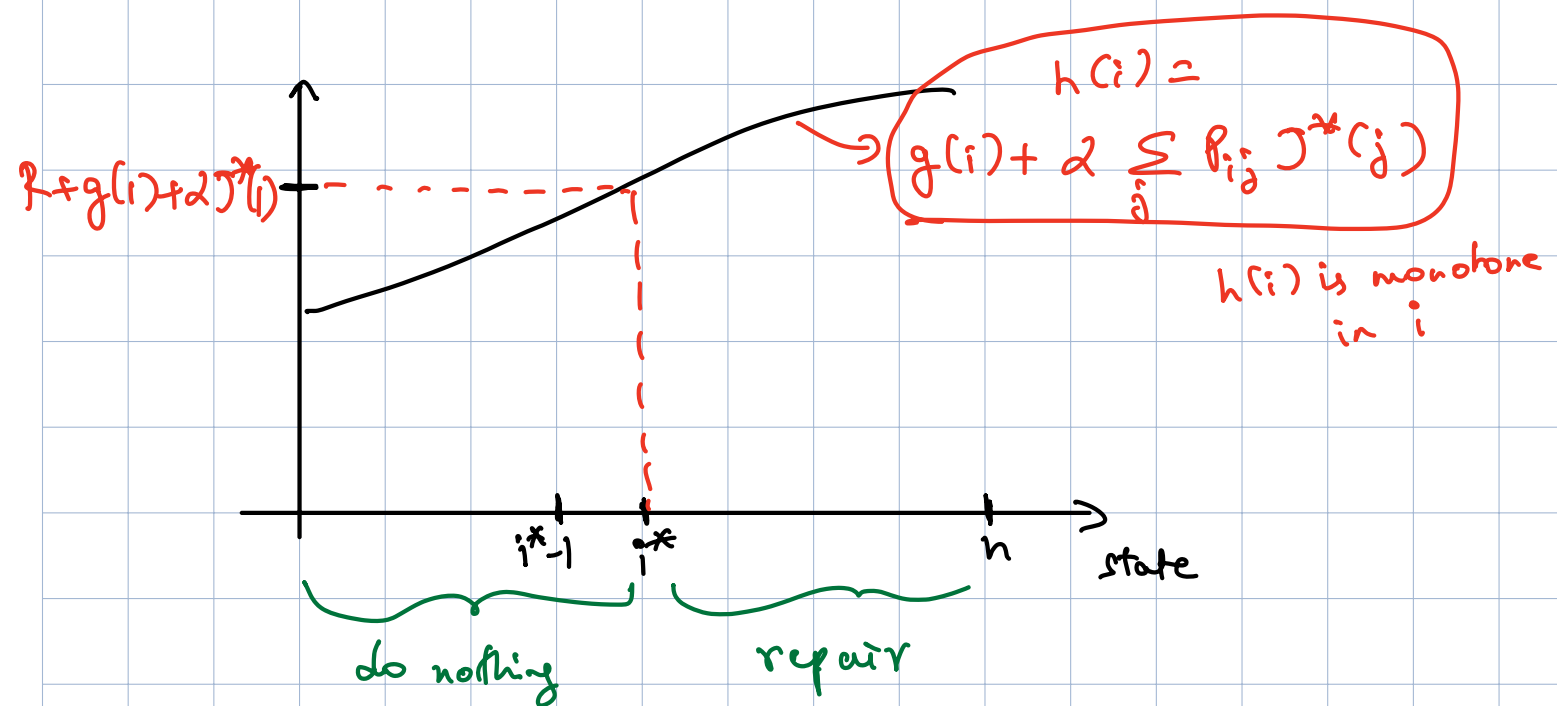
S_R = set of states where it is optimal to repair

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

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 PI converges to optimal policy in at most n iterations.

Illustrative example for vI :

M DP $X = \{1, 2\}$ $A = \{a, b\}$

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

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

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

Discount $\alpha = 0.9$

$$J_0 = (0, 0)$$

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

$$J_1 = TJ_0 = (0.5, 1)$$

$$J_2 = (1.28, 1.56)$$

∴ So on.

PI algorithm:

Step 1: Start with a policy π_0

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

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

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

Step 3: Policy improvement

Find a new policy π_{k+1} by

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

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

If $J_{\pi_{k+1}}(i) < J_{\pi_k}(i)$ 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 policies s.t.

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

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

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:

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

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

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

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

Discount $\alpha = 0.9$

Initialization!

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

Policy evaluation:

Finding J_{π_0} :

Two equations in two unknown $J_{\pi_0}(1)$ & $J_{\pi_0}(2)$

$$J_{\pi_0}(1) = g(1, a) + \alpha P_{11}(a) J_{\pi_0}(1) + \alpha P_{12}(a) J_{\pi_0}(2)$$

$$J_{\pi_0}(2) = g(2, b) + \alpha P_{21}(b) J_{\pi_0}(1) + \alpha P_{22}(b) J_{\pi_0}(2)$$

Using MDP data, we have

$$J_{\pi_0} = T_{\pi_0} J_{\pi_0}$$

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

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

Solving, $J_{\pi_0}(1) = 24.12$, $J_{\pi_0}(2) = 25.96$

Policy improvement!

$$T_{\pi_1} J_{\pi_0} = T J_{\pi_0}$$

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

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

$$\pi_1(1) = b$$

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

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

$$\underline{\underline{\pi_1(2) = a}}$$

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

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

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

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 \rightarrow optimal policy π^*

Optimal cost J^*

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 $J \leq T J$
 Then, $J \leq T^2 J$
 \vdots
 $J \leq T^k J$

$\Rightarrow J \leq J^* = T J^*$ i.e.,

$$\begin{matrix} J(1) \leq J^*(1) \\ \vdots \\ J(n) \leq J^*(n) \end{matrix}$$

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

Optimization problem:

constraints: $J \leq T J$

objective: $\max J$

$$\begin{aligned} x &\leq \min_{y \in \mathcal{Y}_1, \dots, \mathcal{Y}_k} f(y) \\ x &\leq f(y_1) \\ x &\leq f(y_2) \\ &\vdots \\ x &\leq f(y_k) \end{aligned}$$

More precisely, $J \leq T J$
 \Leftrightarrow

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

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

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

$\} \text{ } \left(\begin{array}{l} \text{linear func} \\ \{ \lambda_1, \dots, \lambda_n \} \\ \text{linear functions} \end{array} \right)$

$\} \text{ } \left(\begin{array}{l} \text{constraint} \\ \text{ } \end{array} \right)$

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

variables = n ← 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.