

Infinite-horizon MDPs

Lecture 5* (contd)

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

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

initial state
 ↗
 policy $\pi = \{\mu_0, \mu_1, \mu_2, \dots\}$
 ↗
 discount factor $\alpha \in (0, 1]$
 ↗
 single-stage cost
 "stationary"
 ↗

Goal: $J^*(x_0) = \min_{\pi \in \Pi} J_\pi(x_0)$
 optimal expected cost → set of admissible policies

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

↗
 optimal policy

Two popular MDPs:

(I) Stochastic shortest path (SSP)

(i) "finite state space" $\{1, \dots, n\}$ & finite action space

(ii) $\alpha = 1$

(iii) \exists a state "T" s.t. $p_{TT}(a) = 1$ "absorbing"

$g(T, a, T) = 0 \quad \forall a$
 "cost-free"

(II) Discounted MDP

(i) discount $\alpha < 1$

(ii) $|g(\cdot, \cdot, \cdot)| < M < \infty$

(i) & (ii) $\Rightarrow J_{\pi}(x_0)$ is finite.

(III) Average-cost MDP : skipped. (see Chapter 5 of V.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 & stationary cost $g(i, a, j)$

alert: this is a change of notation from finite horizon chapter

Then, the infinite horizon optimal expected cost can be obtained by

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

(II) Bellman equation

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

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

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

\uparrow
k+1 stage
optimal cost

\uparrow
k-stage
optimal cost

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

$$\textcircled{2} \quad J^*(i) = \min_{a \in A(i)} \sum_{j=1}^n P_{ij}^a(a) (g(i, a, j) + \alpha J^*(j))$$

\hookrightarrow Bellman equation

Eq \textcircled{1} is an algorithm

Eq \textcircled{2} is a system of equations satisfied by J^*

Lecture - 6

\textcircled{3} How to get the optimal policy π^* ?

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

\leftarrow Bellman equation

Make up a policy as follows: For state i , let $a^*(i)$ be the minimizer in the Bellman equation.

Set $\pi(i) = \hat{\alpha}(i)$, $\forall i \in S$.

Then, " $\pi = \pi^*$ " i.e., π is a optimal policy.

Stochastic shortest path (SSP) problems

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

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

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

absorbing cost-free

Examples:

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

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

States are tuples: (i, k)

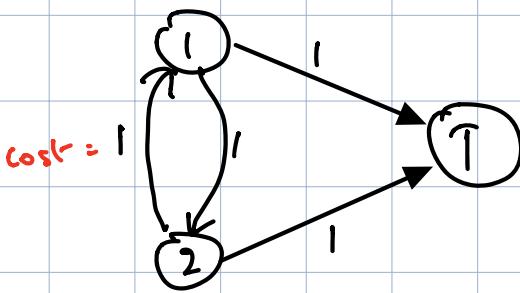
↑ State ↑ Stage

Terminal states: $(i, N) \quad \forall i$

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

Costs: H.W.

Proper policy:



Shortest path problem

Improper policy: loop between 1 & 2

Expected cost = ∞ .

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 π is **proper**

$$\rho_\pi = \max_{i=1, \dots, n} P(i_n \neq T \mid i_0 = i, \pi) < 1$$

↑ ↑ ↑
 nth stage starting state actions chosen
 using π

A policy that ain't proper is "improper" (i.e., $\rho_\pi = 1$)

Assumptions:

(A1) There exists at least one proper policy.

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

Under a proper policy, $\text{Prob}(\text{not reaching } T)$ goes down asymptotically.

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

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

use Markov property & time homogeneity
to conclude this prob $\leq e_\pi^2$

Lecture-7*

More generally,

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

$k < n$ $k \geq n$
 \downarrow \downarrow
 $P(i_k \neq T \mid i_{\lfloor \frac{k}{n} \rfloor} \neq T, i_0 = i, \pi)$
 $\times P(i_{\lfloor \frac{k}{n} \rfloor} \neq T \mid i_0 = i, \pi)$
 $\leq 1 \times e_\pi^{\lfloor \frac{k}{n} \rfloor} = e_\pi^{\lfloor \frac{k}{n} \rfloor}$

Let g_k be the cost incurred in k th stage for policy π

$$E(g_k) \leq e_\pi^{\lfloor \frac{k}{n} \rfloor} \max_{i=1 \dots n} |g(i, \pi(i))|$$

for simplicity,
assume stage
cost is a
function of current
state & action

$$|\mathcal{J}_\pi(i)| \leq \lim_{N \rightarrow \infty} \sum_{k=0}^{N-1} \ell_\pi^{\left[\frac{k}{n}\right]} \max_{j=0..n} \lg(j, \pi(i, j))$$

$$< \infty. \text{ (since } \ell_\pi < 1)$$

So, for a proper policy π , the expected cost \mathcal{J}_π is finite.

Notation :

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

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

as

$$(TJ)(i) = \min_{a \in A(i)} \sum_{j \in S} p_{ij}(a) (g(i, a, j) + J(j)),$$

$i = 1, \dots, n$

↑
includes terminal state

Alternatively,

$$(TJ)(i) = \min_a E_j [g(i, a, j) + J(j)]$$

Another operator for a given policy π :

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

$$T_\pi J(i) = \sum_{j \in S} p_{ij}(\pi(i)) (g(i, \pi(i), j) + J(j)),$$

$i = 1, \dots, n$

For a given π ,

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

P^π matrix has positive entries & row sum ≤ 1

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

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

Using the notation above,

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

Let $(T^k J)(i) = (T(T^{k-1} J))(i)$, $i = 1 \dots n$
 with $(T^0 J)(i) = J(i)$

Similarly, $(T_\pi^k J)(i) = (T_\pi(T_\pi^{k-1} J))(i)$, $i = 1 \dots n$
 with $(T_\pi^0 J)(i) = J(i)$

" 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 $k=1, 2, \dots$

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

(ii) For any stationary policy π ,

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

Pf:

$$\begin{aligned} (T\tau)(i) &= \min_a \sum_{j \in S} p_{ij}(a) (g(i, a, j) + \tau(j)) \\ &\leq \min_a \sum_{j \in S} p_{ij}(a) (g(i, a, j) + \tau'(j)) \\ &= (T\tau')(i). \end{aligned}$$

H.W.

Complete the rest of the proof using induction.

H.W.

Do the proof for T_π .



Another lemma:

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

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

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

Note: Inequalities reversed if $\delta < 0$.

P.F.:

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

$$= \min_a \sum_{j \in X} p_{ij}(a) (g(i, a, j) + (\mathcal{J} + \delta e)(j))$$

$$= \min_a \sum_{j \in X} p_{ij}(a) (g(i, a, j) + \mathcal{J}(j) + \delta)$$

$$= \min_a \left[\sum_{j \in X} p_{ij}(a) (g(i, a, j) + \mathcal{J}(j)) + \sum_j p_{ij}(a) \delta \right]$$

since $\sum p_{ij}(a) \leq 1$

$$\leq \min_a \left[\sum_{j \in X} p_{ij}(a) (g(i, a, j) + \mathcal{J}(j)) + \delta \right]$$

$$= \min_a \left[\sum_{j \in X} p_{ij}(a) (g(i, a, j) + \mathcal{J}(j)) \right] + \delta$$

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

H.W. Complete the proof by induction.

■

Properties of T_π :

Proposition 1:

Assume (A1) & (A2). Then,

(i) For any proper policy π , the associated cost \mathcal{J}_π satisfies

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

for any \mathcal{J} .

$$\text{Also, } J_\pi = T_\pi J_\pi \quad (\infty)$$

& J_π is the unique solution of (∞)

[For computing J_π , one can either apply T_π repeatedly a large number of times (or) solve $J = T_\pi J$]

Policy evaluation: task of computing J_π

(ii) Suppose a stationary policy π satisfies

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

Then, π is proper.

Properties of Bellman optimality operator T :

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

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

π^* \in optimal policy

Then,

(i) J^* satisfies

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

& J^* is the unique solution to (∞)

(ii) For any $\text{finite } J$, we have

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

(iii) A stationary policy π is optimal if & only if

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

Lecture-8*

<Proof of Prop. 1> part(i)

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

$$<\!\!\text{if within Pf}\!\!> T_\pi J = g_\pi + P_\pi J$$

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

Generalize,

$$T_\pi^k J = P_\pi^k J + \sum_{m=0}^{k-1} P_\pi^m g_\pi \quad \rightarrow \infty$$

" π is proper" (Given). So,

$$P(i_k \neq \bar{i} \mid i_0 = i, \pi) \leq e^{-\lfloor \frac{k}{n} \rfloor} \quad \forall i$$

$$\lim_{k \rightarrow \infty} P_\pi^k J = 0$$

So, wifg (*), $\lim_{k \rightarrow \infty} T_\pi^k J = \lim_{k \rightarrow \infty} \sum_{m=0}^{k-1} P_\pi^m g_\pi = J_\pi$ (See defn of J_π below)

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

<End of PF within PF>

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

(PF within PF, again)

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

Taking limit as $k \rightarrow \infty$ on both sides, and using first claim,

$$J_\pi = g_\pi + P_\pi J_\pi$$

$$J_\pi = T_\pi J_\pi$$

Uniqueness: Suppose J satisfies $J = T_\pi J$

$$J = T_\pi^2 J = \dots = T_\pi^k J$$

Taking limit as $k \rightarrow \infty$

$$J = \lim_{k \rightarrow \infty} T_\pi^k J = J_\pi$$

<End of PF within PF, again>

<Proof of Prop 1 - part (ii)>

Given: Stationary π satisfying $J \geq T_\pi J$ for some J

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

$$(or) \quad J \geq T_\pi^2 J \leftarrow \text{so on}$$

$$J \geq T_{\pi}^k J = P_{\pi}^k J + \sum_{m=0}^{k-1} P_{\pi}^m g_{\pi} \quad (\text{from CX a couple of pages ago})$$

If π is not proper, then

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

This leads to a contradiction because

$J(i)$ is finite $\forall i$:

(End of part (ii))

(Proof of Prop 2)

(i) *to show:* J^* satisfies

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

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

(pf) "T has at most one fixed point."

Suppose J and J' satisfy $J = TJ$, $J' = TJ'$

Select a policy π s.t.

$$TJ(i) = T_{\pi} J(i), \quad \forall i$$

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

"Set $\pi(i)$ to be the minimizer on LHS, $\star\star$ ".

Similarly, select a policy π' s.t. $TJ' = T_{\pi'} J'$

Notice that

$$J = T J = T_{\pi} J,$$

So, $J = T_{\pi} J$. Using part(ii) of Prop', we can infer π is proper.

Similarly, $J' = T_{\pi'} J'$ & π' is proper.

Now, $J = T J = T^k J \leq T_{\pi'}^k J$, $\forall k \geq 1$
by definition of T, T_{π} .

Taking limit, $J \leq \lim_{k \rightarrow \infty} T_{\pi'}^k J = J_{\pi'} = J'$

So, $J \leq J'$

Swapping J & J' & repeating the arguments,

$$J' \leq J$$

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

(End of "uniqueness" part)

(Begin of "existence" part)

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

Let π' be another policy s.t.

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

Now,

$$\mathcal{J}_\pi = T_\pi \mathcal{J}_\pi \geq T \mathcal{J}_\pi = T_\pi \mathcal{J}_\pi$$

\mathcal{J}_π is unique
fixed pt of T_π

Def of T, T_π

by construction of π'

So, $\mathcal{J}_\pi \geq T_\pi \mathcal{J}_\pi \Rightarrow \pi'$ is proper
(Prop 1, part (ii))

$$\mathcal{J}_\pi \geq T_\pi^2 \mathcal{J}_\pi \geq \dots \geq T_\pi^k \mathcal{J}_\pi$$

Taking limits,

$$\mathcal{J}_\pi \geq \lim_{k \rightarrow \infty} T_\pi^k \mathcal{J}_\pi = \mathcal{J}_{\pi'}$$

So, we got two proper policies π, π' s.t.

$$\mathcal{J}_\pi \geq \mathcal{J}_{\pi'}$$

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

s.t. (i) Each π_k is proper

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

State-action spaces finite \Rightarrow # policies is finite

\Rightarrow Some policy, say π ,
must be repeated in $\{\pi_k\}$

\Rightarrow w.r.t (ii), we obtain

$$\mathcal{J}_\pi = T \mathcal{J}_\pi$$

< End of existence proof)

Claim: $T^k \mathcal{J} \rightarrow \mathcal{J}_\pi$ as $k \rightarrow \infty$ (Later we show)
 $\mathcal{J}_\pi = \mathcal{J}^*$

CDF: $e = n$ -vector of ones. Fix $\delta > 0$

Pick $\hat{\mathcal{J}}$ s.t. $T_\pi \hat{\mathcal{J}} = \hat{\mathcal{J}} - \delta e$

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

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

i.e., $\hat{\mathcal{J}}$ is the expected cost of a policy in an SSP with single stage cost $(g_\pi + \delta e)$

Since π is proper, $\hat{\mathcal{J}}$ is unique (Prop 1).
part (ii)

By construction $\mathcal{J}_\pi \leq \hat{\mathcal{J}}$ since $g_\pi \leq g_\pi + \delta e$, $\delta > 0$

Lecture 9

Recall \mathcal{J}_π satisfies $\mathcal{J}_\pi = T \mathcal{J}_\pi$

Notice that

$$\mathcal{J}_\pi = T \mathcal{J}_\pi \leq T \hat{\mathcal{J}} \leq T_\pi \hat{\mathcal{J}} = \hat{\mathcal{J}} - \delta e \leq \hat{\mathcal{J}}$$

$$\Rightarrow \mathcal{J}_\pi = T^k \mathcal{J}_\pi \leq T^k \hat{\mathcal{J}} \leq T^{k-1} \hat{\mathcal{J}} \leq \hat{\mathcal{J}}$$

$T \hat{\mathcal{J}} \leq \hat{\mathcal{J}}$
$T^{k-1}(T \hat{\mathcal{J}}) \leq T^{k-1} \hat{\mathcal{J}}$
$T^k \hat{\mathcal{J}} \leq T^k \hat{\mathcal{J}}$

$$\Rightarrow \mathcal{I}_\pi \leq T^k \hat{\mathcal{J}} \leq \hat{\mathcal{J}}$$

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

Hence,

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

Apply T on both sides of the shaded equation

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

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

($T \mathcal{J} = \min_a$ "linear fn(\mathcal{J})" & hence T is continuous)

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

$$\text{So, } T \tilde{\mathcal{J}} = \tilde{\mathcal{J}}$$

$\Rightarrow \tilde{\mathcal{J}} = \mathcal{J}_\pi$ since the fixed point of T is unique & $= \mathcal{J}_\pi$.

"The sandwich principle"

We will show

$$\lim_{k \rightarrow \infty} T^k (\mathcal{J}_\pi - \delta e) = \mathcal{J}_\pi$$

To see this,

See a lemma on p. q above

$$\mathcal{J}_\pi - \delta c = T\mathcal{J}_\pi - \delta e \stackrel{\downarrow}{\leq} T(\mathcal{J}_\pi - \delta c) \leq T\mathcal{J}_\pi = \mathcal{J}_\pi$$

$$\mathcal{J}_\pi - \delta c \leq T(\mathcal{J}_\pi - \delta c) \leq \mathcal{J}_\pi$$

Using " $T(\mathcal{J}_\pi - \delta c) \leq T^2(\mathcal{J}_\pi - \delta c) \leq \dots \leq T^k(\mathcal{J}_\pi - \delta c)$ ",

$$\mathcal{J}_\pi - \delta c \leq T^k(\mathcal{J}_\pi - \delta c) \leq \mathcal{J}_\pi$$

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

$$\lim_{k \rightarrow \infty} T^k(\mathcal{J}_\pi - \delta c) = \mathcal{J}_\pi \quad \xrightarrow{\text{use an argument similar to (x)}}$$

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

$$\mathcal{J}_\pi - \delta c \leq \mathcal{J} \leq \hat{\mathcal{J}} \quad \xrightarrow{\text{cost of policy } \pi \text{ with single stage cost } g_\pi + \delta c}$$

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

Taking limit as $k \rightarrow \infty$,

$$\mathcal{J}_\pi \leq \lim_{k \rightarrow \infty} T^k \mathcal{J} \leq \mathcal{J}_\pi$$

$$\text{So, } \lim_{k \rightarrow \infty} T^k \mathcal{J} = \mathcal{J}_\pi \text{ for any } \mathcal{J}.$$

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

Take any policy π' . We have

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

where J_0 is an arbitrary n -vector.

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

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

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 π is optimal

$$J_\pi = J^*$$

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

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

(\Leftarrow) Suppose $\mathcal{J}^* = T \mathcal{J}^* = T_\pi \mathcal{J}^*$
 Using $\mathcal{J}^* = T_\pi \mathcal{J}^*$, we can infer π is proper
 (Prop 1 part(iii))

For a proper π , if we have

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

$\mathcal{J}^* = \mathcal{J}_\pi$ (Prop 1 part(i))

$\Rightarrow \pi$ 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 $\mathcal{J}^* \rightarrow$ optimal expected cost satisfies

Bellman equation: $\hat{\mathcal{J}}(i) = T \mathcal{J}^*(i)$, which is

$$\mathcal{J}^*(i) = \min_{a \in A(i)} \left[1 + \sum_{j=1}^n p_{ij}(a) \mathcal{J}^*(j) \right], \forall i$$

Special Case: Only one action in each state \leftarrow Markov chain

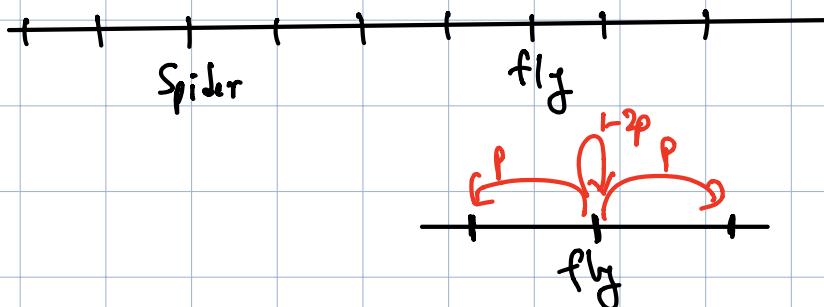
$T^*(i) \leftarrow$ expected first passage time to terminal state "T".

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

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

Another example: Fly & a spider

Spider & fly move on a line.



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

(ii) if distance $= 1$

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 = $\{0, 1, \dots, n\}$

\uparrow
initial distance

Terminal state = 0

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

\uparrow
spider moves

\downarrow
spider doesn't move

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

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

For $i \geq 2$,

$$P_{ii} = p, \quad P_{i,i-1} = 1-2p, \quad P_{i,i-2} = p$$

In state 1,

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

Spider doesn't jump: $P_{12}(\bar{M}) = p, \quad P_{11}(\bar{M}) = 1-2p, \quad P_{10}(\bar{M}) = p$

Bellman equation for $i \geq 2$

$$\mathcal{J}^*(i) = 1 + p \mathcal{J}^*(i) + (1-2p) \mathcal{J}^*(i-1) + p \mathcal{J}^*(i-2)$$

$$\mathcal{J}^*(0) = 0$$

\square

For state = 1 :

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

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

From (\star) ,

$$\begin{aligned} \mathcal{J}^*(2) &= 1 + p \mathcal{J}^*(2) + (1-2p) \mathcal{J}^*(1) \\ \Rightarrow \mathcal{J}^*(2) &= \frac{1}{1-p} + \frac{(1-2p) \mathcal{J}^*(1)}{1-p} - (\star\star\star) \end{aligned}$$

Substitute $(\star\star\star)$ in (\star) ,

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

which is the same as

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

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

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

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

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

Substitute this in (****) to obtain optimal to jump.

$$\frac{2p}{1-2p} \leq \frac{P}{1-p} + \frac{1}{1-p} \Leftrightarrow p \leq \frac{1}{3}$$

(case (A) > (B))

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

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

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

Substitute this in the constraint to obtain

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

optimal to
not jump.

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

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

Coming next!

Value iteration!

Start with π_0

$$\pi_0 \xrightarrow{T} \pi_1 \xrightarrow{T} \pi_2 \dots$$
$$\pi_{k+1} = T \pi_k$$

We know under (A1) & (A2),
(at least one proper policy)
(improper policies have
no cost)

$$\lim_{k \rightarrow \infty} T^k \pi_0 = \pi^*$$

Special Case: "All policies are proper"

Then, VI converges at a geometric rate.

$$\|\pi_k - \pi^*\|_\infty \leq e^k \|\pi_0 - \pi^*\|_\infty$$

$0 < e < 1$

← need theory of
contraction
mappings.

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. $\ell(i) > 0 \forall i=1..n$, and a scalar $\rho < 1$

such that

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

modulus of contraction

Here

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

Weighted max-norm

Also, for any stationary proper policy π ,

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

Pf:

Need: a vector ξ such that T is a contraction wrt $\|\cdot\|_{\xi}$

Idea: Make up a MDP, solve it & get ξ .

& then show such a ξ helps in contraction business.

Consider a new SSP where transition probabilities

are the same as in the original SSP, but

transition costs = -1 except at terminal state.

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

Let $\hat{\pi}_j \rightarrow$ optimal cost in this new SSP.

Thus, using Bellman equation

$$\hat{\pi}_i = -1 + \min_{a \in A(i)} \sum_{j=1}^n p_{ij}(a) \hat{\pi}_j$$

$$= T \hat{\pi}_i$$

$$\leq T_\pi \hat{\pi}_i \quad \text{for any } \pi \text{ (proper stationary)}$$

$$= -1 + \sum_{j=1}^n p_{ij}(\pi(i)) \hat{\pi}_j$$

$$\text{Let } \xi(i) = -\hat{\pi}_i, \quad i = 1 \dots n$$

$$\text{Note: } ① \quad \xi(i) \geq 1 \quad \forall i$$

$$② \quad -\hat{\pi}_i \geq 1 + \sum_{j=1}^n p_{ij}(\pi(i)) (-\hat{\pi}_j)$$

$$\Rightarrow \xi(i) \geq 1 + \sum_{j=1}^n p_{ij}(\pi(i)) \xi(j)$$

$$\sum_{j=1}^n p_{ij}(\pi(i)) \xi(j) \leq \xi(i) - 1 \leq \rho \xi(i), \quad \text{---(x)}$$

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

Now, for any $\pi, \tau, \bar{\tau}$, $T_\pi \tau(i) = \sum_j p_{ij}(\pi(i)) \xi(j) + \tau(j)$

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

$$\leq \sum_{j=1}^n p_{ij}(\pi(i)) |\tau(j) - \bar{\tau}(j)| = \sum_{j=1}^n p_{ij}(\pi(i)) \xi(j) \frac{|\tau(j) - \bar{\tau}(j)|}{\xi(j)}$$

$$\leq \left(\sum_{j=1}^n p_{ij}(\pi(i)) \xi(j) \right) \left(\max_{j=1 \dots n} \frac{|\tau(j) - \bar{\tau}(j)|}{\xi(j)} \right)$$

\checkmark using (x)

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

holds for any i

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

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

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

So, T_π is a contraction w.r.t. $\|\cdot\|_\xi$

with modulus e .

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

We have shown already that

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

Take minimum over π on both sides to obtain

$$(T \mathcal{J})(i) \leq (T \bar{\mathcal{J}})(i) + C \xi(i) \max_{j=1 \dots n} \frac{|\mathcal{J}(j) - \bar{\mathcal{J}}(j)|}{\xi(j)}$$

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

$$(T \bar{\mathcal{J}})(i) \leq (T \mathcal{J})(i) + C \xi(i) \max_{j=1 \dots n} \frac{|\mathcal{J}(j) - \bar{\mathcal{J}}(j)|}{\xi(j)}$$

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

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

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

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

Value iteration (VI):

1) Choose some \mathcal{J}_0

2) Repeatedly apply T

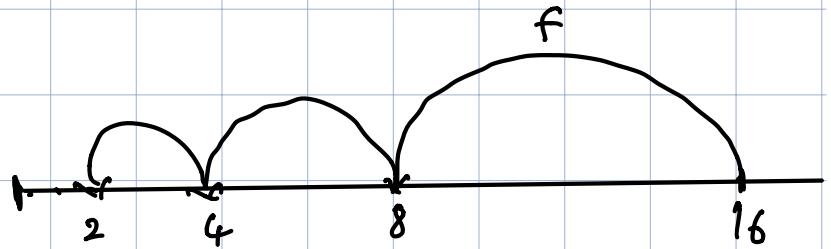
$$\mathcal{J}_0 \xrightarrow{T} \mathcal{J}_1 \xrightarrow{T} \mathcal{J}_2 \xrightarrow{T} \dots$$

$$\mathcal{J}_{k+1} = T \mathcal{J}_k$$

We know $\lim_{k \rightarrow \infty} T^k \mathcal{J}_0 = \mathcal{J}^*$ for any \mathcal{J}_0 .

Contraction mappings: A quick detour

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



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

$$\underbrace{f(0) = 0}$$

zero is a fixed point

& I get to the fixed point by repeatedly applying f .

Vector space X .

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

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

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

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

Contraction mapping!

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

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

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

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 $\exists N$ s.t.

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

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

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

(ii) $x_{k+1} = F(x_k)$. Then

$$x_k \rightarrow x^* \text{ as } k \rightarrow \infty.$$

Lecture-12

Contraction mappings in MDPs: (Ref: Sec 1.5 of DPoC-VcII).

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

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

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

Simple case: finite state space.

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

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

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

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

$$\mathcal{T}^* = F \mathcal{T}^*$$

(ii) $\lim_{k \rightarrow \infty} F^k \mathcal{T}_0 = \mathcal{T}^*$ VI converges

$$\text{At } k_0, \quad \|F^{k_0} \mathcal{T}_0 - \mathcal{T}^*\|_{\varsigma} \leq e^k \|\mathcal{T}_0 - \mathcal{T}^*\|_{\varsigma}.$$

error bound for VI

Pf:

See next page

Fix some $J_0 \in B(\mathcal{X})$.

Do $J_{k+1} = F J_k$ starting with J_0

$$\|J_{k+1} - J_k\|_{\xi} \leq \ell \|J_k - J_{k-1}\|_{\xi}$$

F is a ℓ -contraction

$$\Rightarrow \|J_{k+1} - J_k\|_{\xi} \leq \ell^k \|J_1 - J_0\|_{\xi} \quad -(*)$$

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

$$\|J_{k+m} - J_k\|_{\xi} = \|(J_{k+m} - J_{k+m-1}) + (J_{k+m-1} - J_{k+m-2}) + \dots + (J_{k+1} - J_k)\|_{\xi}$$

triangle
 $i \neq j$

$$\leq \sum_{i=1}^m \|J_{k+i} - J_{k+i-1}\|_{\xi}$$

using (*)

$$\leq \ell^k (1 + \ell + \dots + \ell^{m-1}) \|J_1 - J_0\|_{\xi}$$

$$\|J_{k+m} - J_k\|_{\xi} \leq \frac{\ell^k}{1-\ell} \|J_1 - J_0\|_{\xi}$$

Is $\{J_k\}$ a Cauchy sequence? Yes.

Since $B(\mathcal{X})$ is complete, $J_k \rightarrow J^*$ and $J^* \in B(\mathcal{X})$.

To show: \mathcal{J}^* is a fixed point.

$$\|F\mathcal{J}^* - \mathcal{J}^*\|_{\xi} \leq \|F\mathcal{J}^* - \mathcal{J}_k\|_{\xi} + \|\mathcal{J}_k - \mathcal{J}^*\|_{\xi}$$

using $\mathcal{J}_k = F\mathcal{J}_{k-1} \rightarrow$
& F is contraction $\leq \rho \|\mathcal{J}^* - \mathcal{J}_{k-1}\|_{\xi} + \|\mathcal{J}_k - \mathcal{J}^*\|_{\xi}$

$\rightarrow 0$ as $k \rightarrow \infty$ since $\mathcal{J}_k \rightarrow \mathcal{J}^*$.

So, $F\mathcal{J}^* = \mathcal{J}^*$.

To show: uniqueness of \mathcal{J}^* .

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

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

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

$$\Rightarrow \mathcal{J}^* = \tilde{\mathcal{J}} \text{ (why? because } \rho < 1)$$

To show: error bound

$$\|F^k \mathcal{J}_0 - \mathcal{J}^*\|_{\xi} = \|F^k \mathcal{J}_0 - F\mathcal{J}^*\|_{\xi}$$

$$\leq \rho \|\mathcal{J}_0 - \mathcal{J}^*\|_{\xi}$$

& repeat to infer

$$\|F^k \mathcal{J}_0 - \mathcal{J}^*\|_{\xi} \leq \rho^k \|\mathcal{J}_0 - \mathcal{J}^*\|_{\xi}$$

■

Back to SSPs:

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

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

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

$$\| \pi_k - \pi^* \|_{\xi} \leq \rho^k \| \pi_0 - \pi^* \|_{\xi},$$

where $\pi_k = T\pi_{k-1}$ & π_0 is the initial vector
for value iteration.

Remark: A similar error bound holds for evaluating
a proper policy π using VI.

Gauss-Seidel variant of VI:

Note: In VI, we do $\pi_{k+1} = T\pi_k$, 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:

$$(F\mathcal{J})(_1) = \min_{a \in f(1)} \sum_j p_{1,j}(a) (g(1,a,j) + \mathcal{J}(j))$$

Same as $(T\mathcal{J})(1)$

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

$$(F\mathcal{J})(i) = \min_{a \in f(i)} \left[\sum_j p_{i,j}(a) g(i,a,j) + \sum_{j=1}^{i-1} p_{i,j}(a) (F\mathcal{J})(j) + \sum_{j=i+1}^n p_{i,j}(a) \mathcal{J}(j) \right]$$

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

Remark: Gauss-Seidel VI converges to \mathcal{J}^* faster than regular VI. See 2.2.2 of DPOC-Vol II.

$$\|\mathcal{J}_k^{VI} - \mathcal{J}^*\|_1 \geq \|\mathcal{J}_k^{GSrI} - \mathcal{J}^*\|_1 \leftarrow \text{See Prop 2.2.3 here.}$$

Asynchronous VI

(i) Start with an arbitrary \mathcal{J}_0

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

$$\mathcal{J}_{k+1}(i) = \begin{cases} (T\mathcal{J}_k)(i) & \text{if } i = i_k \\ \mathcal{J}_k(i) & \text{else} \end{cases}$$

Assuming all states in $\{1 \dots n\}$ are picked infinitely often,
it can be shown that

$$\mathcal{T}_k \rightarrow \mathcal{T}^* \text{ as } k \rightarrow \infty.$$

(See the same section of VI)

Lecture 13*

Q-learning as a form of VI

(Visual Q-learning doesn't require the model. But we are following DQC Vol II nomenclature)

When the model is known, i.e., we can compute T ,

Q-learning is equivalent to VI.

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

$$VI : \mathcal{T}_{k+1} = T \mathcal{T}_k$$

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

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

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

Using Bellman equation $J = T J$:

$$J^*(i) = \min_{a \in A(i)} Q^*(i, a)$$

Then, Q^* can be alternatively defined by

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

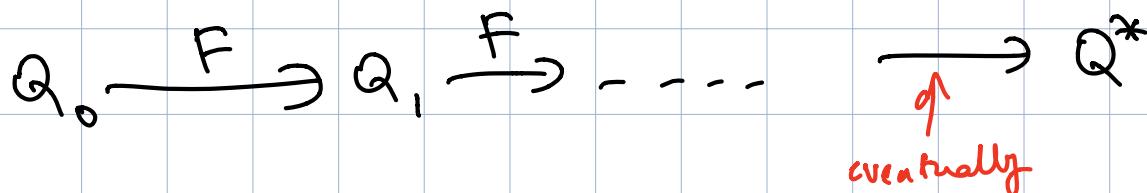
Q-Bellman equation

The operator underlying Q-Bellman equation is

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

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



Once you have Q^* , we can obtain π^* by

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

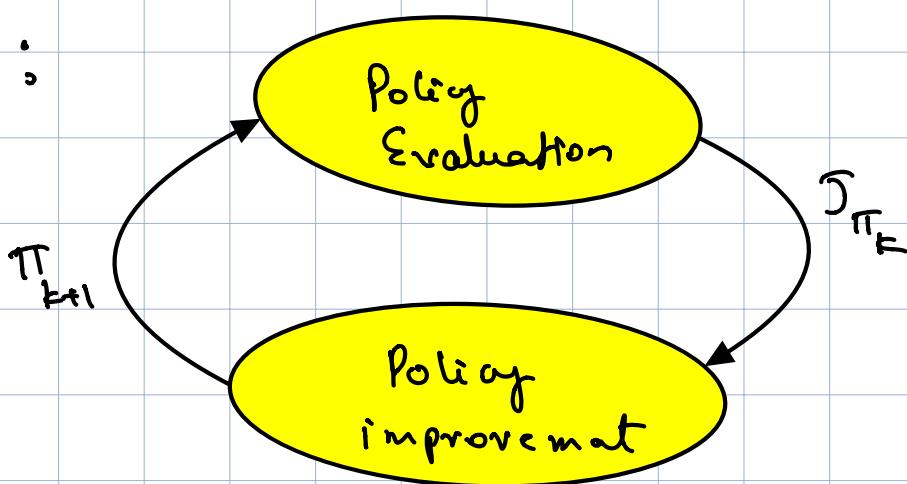
Eq (1-3) is just VI on Q-values.

Policy Iteration (PI)

VI can possibly take infinite # of iterations to converge

PI is a method "guaranteed" to converge within a finite # of iterations.

PI :



$$J_{\pi_0} \geq J_{\pi_1} \dots$$

↑
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.

PI algorithm:

Step 1: Start with a proper policy π_0

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

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

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

Step 3: Policy improvement

Find a new policy π_{k+1} by

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

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

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

then go to Step 2 & repeat.

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

Step 1 can be modified by starting with a finite \mathcal{J} & going to Step 3 directly. That would give π_1 , & from thereon do Step 2 & 3 in tandem until convergence.

Claim: Policy improvement does just that.

Let π, π' be two proper 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:

We know that $J_{\pi} = T_{\pi} J_{\pi}$

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

So,

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

$$\begin{aligned} & \stackrel{\text{because } T_{\pi'} \geq T \pi}{\geq} \sum_j p_{ij}(\pi'(i)) (g(i, \pi'(i), j) + J_{\pi}(j)) \\ &= (T_{\pi'} J_{\pi})(i) \end{aligned}$$

$$\text{So, } J_{\pi} \geq T_{\pi'} J_{\pi} \dots \geq T_{\pi'}^k J_{\pi} \dots \stackrel{k \rightarrow \infty}{\geq} \lim J_{\pi'}^k J_{\pi} = J_{\pi}'$$

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

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

$$\mathcal{J}_\pi = T_{\pi'}, \mathcal{J}_\pi \quad \text{--- } ①$$

We also have $T_\pi, \mathcal{J}_\pi = T \mathcal{J}_\pi \quad \text{--- } ②$

Combining ① & ②,

$$\mathcal{J}_\pi = T \mathcal{J}_\pi$$

So, π is optimal.

Thus, if π isn't optimal, then $\mathcal{J}_{\pi'}(i) < \mathcal{J}_\pi(i)$
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 generates a better policy
if the latter isn't optimal.

Modified PI:

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

Let π_1, π_2, \dots and π_0, π_1, \dots be computed as follows:

$$\xrightarrow{\text{Policy improvement}} T_{\pi_k} \pi_k = T \pi_k$$

$$\xrightarrow{\text{Policy evaluation}} \pi_{k+1} = T_{\pi_k}^{m_k} \pi_k$$

by m_k steps
of VI

Two Special Cases:

If $m_k = \infty$, Modified PI = PI since $T_{\pi_k}^{\infty} \pi_k = \pi_k$

If $m_k = 1$, Modified PI = VI

$$\text{since } \pi_{k+1} = T_{\pi_k}^{m_k} \pi_k = T_{\pi_k} \pi_k = T \pi_k$$

Recommended: Choose $m_k > 1$

Convergence: Modified PI converges

"proof skipped" See Ch. 2 of

Bertsekas DPOC-Vol II

Lecture - 14*

Asynchronous PI :

Let $\mathcal{J}_1, \mathcal{J}_2, \mathcal{J}_3, \dots$: Sequence of optimal cost estimates
 $\pi_1, \pi_2, \pi_3, \dots$: Corresponding sequence of policies

Given (\mathcal{J}_k, π_k) , we select a subset S_k of states
and generate $(\mathcal{J}_{k+1}, \pi_{k+1})$ in one of the
following two ways:

(Way 1)

$$\mathcal{J}_{k+1}(i) = \begin{cases} (\mathcal{J}_{\pi_k} \mathcal{J}_k)(i) & \text{if } i \in S_k \\ \mathcal{J}_k(i) & \text{else} \end{cases}$$

Set $\pi_{k+1} = \pi_k$

(Way 2)

$$\pi_{k+1}(i) = \begin{cases} \underset{a \in A(i)}{\operatorname{arg\min}} \sum_j p_{ij}^{(a)} (j(i, a, j) + \mathcal{J}_k(j)) & \text{if } i \in S_k \\ \pi_k(i) & \text{else} \end{cases}$$

Set $\mathcal{J}_{k+1} = \mathcal{J}_k$

Special cases of Async-PI:

① If $S_k = \text{entire state space}$, ($*$) is performed infinite # of times before one (+) operation, then we get **regular PI**

② If $S_k = \text{entire state space}$, ($*$) performed m_k times before (+), then we get **modified PI**

③ If $S_k = \text{entire state space}$, one ($*$) operation followed immediately by one (+) operation, then we get **Value iteration**

④ If $|S_{\text{fd}}| = 1$, and one ($*$) operation followed immediately by one (+) operation, then we get **Asynchronous VI**

Remark: Async-PI 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))
PI leads to "Actor-critic methods".