



Notes

Dynamic

Bh-2

★ SSPP :-

- * There exists a cost free transition state 0.

$$p_{00}(u) = 1, \quad g(0, u, 0) = 0 \quad \forall u \in A(0).$$

$$\gamma = 1. \quad (\text{No discounting}).$$

Problem :- How to reach terminal state with minimum expected cost?

Non-terminal states :- 1, 2, ..., n.

Let

$$J_\mu(i) = \lim_{N \rightarrow \infty} E_\mu \left[\sum_{m=0}^{N-1} g(S_m, \mu(S_m), S_{m+1}) \mid S_0 = i \right]$$

Here μ is a stationary policy.

- could be deterministic or randomized.

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

$$\mu_k(i) \in A(i) \quad \forall k, \forall i.$$

* In a stationary policy,

$$\mu_0 = \mu_1 = \mu_2 = \dots = \mu.$$

$$\hookrightarrow J_\mu(i) = \lim_{N \rightarrow \infty} E_\mu \left[\sum_{m=0}^{N-1} g(S_m, \mu(S_m), S_{m+1}) \mid S_0 = i \right].$$

Since costs are bounded :-

$$\max_{i,j} |g(i, \mu(i), j)| \leq K < \infty$$

* Definitions :-

(i) A stationary policy μ is said to be proper if

$$P_\mu = \max_{i=1, \dots, n} P(S_n \neq 0 | S_0 = i, \mu) < 1$$

(ii) A stationary policy that is not proper is called improper.

Note:- μ is proper \Leftrightarrow In the Markov chain corresponding to μ , there exists a positive prob. path from each state to terminal state.

$$\begin{aligned} P(S_n \neq 0 | S_0 = i, \mu) &= P(S_{2n} \neq 0 | S_n \neq 0, S_0 = i, \mu) * \\ &\quad \underbrace{P(S_n \neq 0 | S_0 = i, \mu)}_{\leq P_\mu} \\ &\leq P_\mu. \\ &\leq P_\mu^2 \end{aligned}$$

More generally,

$$P(S_k \neq 0 | S_0 = i, \mu) \leq P_\mu \quad [K/n] \quad i = 1, \dots, n.$$

Suppose:- $n < k < 2n$.

$$P(S_k \neq 0 | S_0 = i, \mu).$$

$$\begin{aligned} &= P(S_k \neq 0 | S_{2n} \neq 0, S_0 = i, \mu) * P(S_n \neq 0 | S_0 = i, \mu) \\ &\leq 1 \leq P_\mu. \end{aligned}$$

Prob. of reaching terminal state after k states $\rightarrow 1$
as $k \rightarrow \infty$ regardless of the initial state

Recall that :-

$$J_\mu(i) = \sum_{m=0}^{\infty} E_\mu [g(s_m, \mu(s_m), s_{m+1}) \mid s_0 = i]$$

$$|J_\mu(i)| \leq \sum_{m=0}^{\infty} E_\mu [|g(s_m, \mu(s_m), s_{m+1})| \mid s_0 = i]$$

$$= \sum_{m=0}^N \sum_j \sum_k p_{ij}(\mu(i)) \cdot p_{jk}(\mu(j)) |g(j, \mu(j), k)|$$

If $\sum_k p_{jk}(\mu(j)) |g(j, \mu(j), k)| = |g_\mu(j)|$

$$= \sum_{m=0}^{\infty} \sum_j p_{ij}(\mu(i)) |g_\mu(j)|$$

$$(g_{\mu}(i))^2 + (g_{\mu}(i))^2 = (g_{\mu}(i))^2$$

$$(g_{\mu}(i))^2 + (g_{\mu}(i))^2 = (g_{\mu}(i))^2$$

$$\frac{(g_{\mu}(i))^2}{(g_{\mu}(i))^2} + \frac{(g_{\mu}(i))^2}{(g_{\mu}(i))^2} = 1$$

$$(g_{\mu}(i))^2 + (g_{\mu}(i))^2 = (g_{\mu}(i))^2$$

* SSPP or Episodic problems :-

$$S = NT \cup T$$

NT = set of non-terminal states.

T = --- terminal states.

$$NT = \{1, 2, \dots, n\}, T = \{0\}$$

proper policy :- μ is proper if

$$P_\mu = \max_{i \in NT} P \{ S_n \neq 0 \mid S_0 = i\} < 1$$

Let $J = (J(0), J(1), \dots, J(n))$, where $J(0) = 0$.

$$\begin{aligned} (TJ)(i) &= \min_{\mu \in A(i)} \sum_{j \in S} p_{ij}(\mu) (g(i, u, j) + J(j)) \\ &= \min_{\mu \in A(i)} \left(p_{i0}(u) g(i, u, 0) + \sum_{j=1}^n p_{ij}^*(u) (g(i, u, j) + J(j)) \right). \end{aligned}$$

Suppose $\bar{g}(i, u) = \sum_{j \in S} p_{ij}(u) \cdot g(i, u, j)$; $\min_{\mu \in A(i)} \left(\bar{g}(i, u) + \sum_{j=1}^n p_{ij}^*(u) J(j) \right)$

$$TJ(i) = \min_{u \in A(i)} \left(\bar{g}(i, u) + \sum_{j=1}^n p_{ij}(u) J(j) \right)$$

Likewise define an operator $T_\mu : \mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1}$ as follows :-

$$(T_\mu J)(i) = \left(\bar{g}(i, \mu(i)) + \sum_{j=1}^n p_{ij}(\mu(i)) J(j) \right)$$

let $P_\mu = \begin{bmatrix} p_{1,1}(\mu(1)) & p_{1,2}(\mu(1)) & \dots & p_{1,n}(\mu(1)) \\ p_{2,1}(\mu(2)) & p_{2,2}(\mu(2)) & \dots & p_{2,n}(\mu(2)) \\ \vdots & & & \\ p_{n,1}(\mu(n)) & p_{n,2}(\mu(n)) & \dots & p_{n,n}(\mu(n)) \end{bmatrix}$

Note :-

$$\sum_{j=1}^n p_{ij}(\mu(p)) \leq 1 \quad \forall i.$$

$$\text{Let } \bar{g}_\mu = (\bar{g}(i, \mu(i)), i \in NT)$$

$$\text{Then, } T\mu J = \bar{g}_\mu + P_\mu J.$$

$$\text{Let } (T^K J)(i) = T(T^{K-1} J)(i).$$

$$(T^2 J)(i) = T(TJ)(i)$$

$$= \min_{\mu \in A(i)} \left(\bar{g}(i, u) + \sum_{j=1}^n p_{ij}(u) (TJ)(j) \right).$$

$$= \min_{\mu \in A(i)} \left(\bar{g}(i, u) + \sum_{j=1}^n p_{ij}(u) \left(\min_{v \in A(j)} \left(\bar{g}(j, v) + \sum_{k=1}^n p_{jk}(v) J(k) \right) \right) \right).$$

Note :- $(T^K J)(i)$ = Optimal cost-to-go for a K -stage finite horizon problem with initial state i , one-stage cost \bar{g} and terminal cost J .

* Monotonicity Lemma :-

For any $J, \bar{J} \in \mathbb{R}^{n+1}$ s.t.

$$J(i) \leq \bar{J}(i) \quad \forall i = 1, \dots, n.$$

$$\& J(0) = \bar{J}(0) = 0.$$

and any stationary μ . \rightarrow what is stationary μ ?
same μ for both J, \bar{J} .

$$(i) \quad (T^K J)(i) \leq (T^K \bar{J})(i) \quad \forall 1 \leq i \leq n.$$

$$(ii) \quad (T_\mu^K J)(i) \leq (T_\mu^K \bar{J})(i) \quad \forall 1 \leq i \leq n$$

Proof :-

consider $k=1$,

$$\begin{aligned}(TJ)(i) &= \min_{\mu \in A(i)} \left(\bar{g}(i, \mu) + \sum_{j=1}^n p_{ij}(\mu) \cdot J(j) \right) \\ &\leq \min_{\mu \in A(i)} \left\{ g(i, \mu) + \sum_{j=1}^n p_{ij}(\mu) \bar{J}(j) \right\} \leq (\bar{T}\bar{J})(i)\end{aligned}$$

Assume above holds for $k=k$.

$$\Rightarrow (T^k J)(i) \leq (T^k \bar{J})(i) \quad \forall i=1, 2, \dots, n$$

For $k=k+1$,

$$\begin{aligned}(T^{k+1} J)(i) &= \min_{\mu \in A(i)} \left(\bar{g}(i, \mu) + \sum_{j=1}^n p_{ij}(\mu) (T^k J)(j) \right) \\ &\leq (T^{k+1} \bar{J})(i) \quad \forall i=1, \dots, n.\end{aligned}$$

By induction the claim follows.

Lemma :- (I)

$\forall k \geq 0$ Vectors J , stationary μ . and $r \geq 0$,

$$(i) T^k (J+r e)(i) \leq (T^k J)(i) + r, \quad i=1, \dots, n$$

$$(ii) T_\mu^k (J+r e)(i) \leq (T_\mu^k J)(i) + r, \quad i=1, \dots, n.$$

$$\begin{aligned}J &= (J(i), i=0, 1, \dots, n) \\ J(0) &= 0.\end{aligned}$$

The inequalities are reversed for $r < 0$.

Here $e = (1, 1, \dots, 1)^T$ is the $(n+1)$ vector of all 1 's.

Note that :-

$$\begin{aligned}
 T(J+r\epsilon)(i) &= \min_{\mu \in A(i)} \left(\bar{g}(i, \mu) + \sum_{j=1}^n p_{ij}(\mu) (J+r\epsilon)(j) \right) \\
 &= \min_{\mu \in A(i)} \bar{g}(i, \mu) + \sum_{j=1}^n p_{ij}(\mu) J(j) + \\
 &\quad \underbrace{\sum_{j=1}^n p_{ij}(\mu) r\epsilon(j)}_{\leq 8} \\
 &\leq (TJ)(i) + r.
 \end{aligned}$$

Rest can be shown by induction.

- Assumptions:-
- (A) \exists at least 1 proper policy.
 - (B) \nexists improper μ , $J^K(i) = \infty$ for at least one state $i \in NT$.

Aside :-

$$J_\mu(i) = E_\mu \left[\sum_{k=0}^{\infty} g(x_k, \mu(x_k), x_{k+1}) \mid x_0 = i \right]$$

$$J^*(i) = \min_{\mu} J_\mu(i) \quad \forall i = 1, \dots, n.$$

we shall show that $T^K J \rightarrow J^*$

$$T_\mu^K J \rightarrow J_\mu \text{ as } k \rightarrow \infty.$$

$$(TJ)(i) = \min_{\mu \in A(i)} \left(\bar{g}(i, \mu) + \sum_{j=1}^n p_{ij}(\mu) J(j) \right)$$

$$J^*(i) = (TJ^*)(i) = \min_{\mu \in A(i)} \left(\bar{g}(i, \mu) + \sum_{j=1}^n p_{ij}(\mu) J^*(j) \right)$$

↓
we shall show this

Bellman Equation for Optimality

Similarly,
 $J_\mu(i) = (T_\mu J_\mu)(i)$

Bellman eqn for μ .

★ Proposition 1:

(a) For a proper policy μ , the associated cost vector J_μ satisfies $(T_\mu^K J)(i) \rightarrow J_\mu(i)$ as $K \rightarrow \infty$ & $i = 1, \dots, n$

for every vector J also,

$J_\mu = T_\mu J_\mu$ and J_μ is the unique solⁿ to this equation.

(b) A stationary policy μ satisfying for some vector J ,

$J(i) \geq (T_\mu J)(i) \quad \forall i = 1, \dots, n$. is proper.

Proof :-

(a) Recall that :-

$$T_\mu J = g_\mu + P_\mu J$$

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

$$= g_\mu + P_\mu (g_\mu + P_\mu J)$$

$$\Rightarrow T_\mu^2 J = g_\mu + P_\mu g_\mu + P_\mu^2 J$$

Proceeding similarly :-

$$T_\mu^K J = P_\mu^K J + \sum_{m=0}^{K-1} P_\mu^m g_\mu.$$

$$\text{Note :- } (P_\mu^K J)(i) = \sum_{j=1}^n P(X_K=j | x_0=i, \mu) J(j)$$

$$\leq \sum_{j=1}^n P(X_K=j | x_0=i, \mu) \left(\max_{j=1, \dots, n} J(j) \right)$$

$$= P(X_K \neq 0 | x_0=i, \mu) \max_{j=1, \dots, n} J(j).$$

$$\leq P_\mu^{(K-1)} \cdot \max_{j=1, \dots, n} J(j)$$

Thus as $K \rightarrow \infty$ the term $\rightarrow 0$ $\therefore p_\mu \leq 1$

$$\therefore \lim_{K \rightarrow \infty} [T_\mu^K J] = \lim_{K \rightarrow \infty} \left\{ \sum_{m=1}^{K+1} P_\mu^m g_\mu \right\} = J_\mu.$$

By definition:

$$T_\mu^{K+1} J = g_\mu + P_\mu (T_\mu^K J)$$

Letting $K \rightarrow \infty$ on either side

$$J_\mu = g_\mu + P_\mu J_\mu \quad \text{or} \quad J_\mu = T_\mu J_\mu.$$

Proof of Uniqueness
By contradiction

Suppose \bar{J} is one more fixed pt. of T_μ .

$$\Rightarrow T_\mu \bar{J} = \bar{J}$$

$$\text{Thus } J_\mu = \bar{J}$$

$$\Rightarrow T_\mu^2 \bar{J} = T_\mu \bar{J} = \bar{J}$$

(b) Given that for a stationary μ & same J .

$$T_\mu^n \bar{J} = \bar{J}$$

$$J(i) \geq (T_\mu J)(i) \quad \forall i = 1, \dots, n.$$

$$\lim_{n \rightarrow \infty} (T_\mu^n \bar{J}) = \bar{J}$$

$$\text{or } J \geq T_\mu J$$

Applying T_μ repeatedly on either side

$$J_\mu$$

$$T_\mu J \geq T_\mu^2 J \geq T_\mu^3 J \geq \dots \geq T_\mu^K J.$$

$$= P_\mu^K J + \sum_{m=0}^{K-1} P_\mu^m g_\mu \rightarrow J_\mu.$$

$$\Rightarrow J(i) \geq J_\mu(i) + i.$$

Now if μ is not proper $J_\mu(i) = \infty$ for at least one $i \Rightarrow J(i) = \infty$ for that!

Converse happens since each $J(i) \in \mathbb{R} \Rightarrow \mu$ is proper.

Proposition 3 :-

(a) The optimal cost vector J^* satisfies the Bellman Equation

$J^* = TJ^*$ & moreover J^* is the unique solⁿ to the eqn.

(b) We have $\lim_{K \rightarrow \infty} (T^K J)(i) = J^*(i) \quad \forall i = 1, \dots, n$.
for every vector J .

(c) A stationary play μ is optimal if & only if $T\mu J^* = T J^*$

$$T = T^* \text{ and } T = \bar{T}$$

$$\bar{T} = \bar{T}^* \text{ and } \bar{T} = \bar{T}^*$$

$$\bar{T} = \bar{T}^* \text{ and } \bar{T} = \bar{T}^*$$

$$\bar{T} = (\bar{T}^*)^*$$

$$T^* \leq T \text{ and } T \leq T^*$$

$$T^* \leq T \leq T^*$$

$$T^* \leq T \leq T^*$$

$$T^* \leq T \leq T^*$$

Lec-9*Proposition : 2

(a) The optimal cost vector J^* satisfies the Bellman eqn $TJ^* = TJ^*$. Moreover, J^* is the unique solⁿ to this eqn.

(b) We have $\lim_{K \rightarrow \infty} (T^K J)(i) = J^*(i) \quad \forall i = 1, \dots, n$

for every vector J .

(c) A stationary policy μ is optimal if and only if $T_\mu J^* = TJ^*$

Proof (a), (b) :-

We first show that T has at most one fixed point.

Suppose there are two fixed points J and J' and suppose μ and μ' are such that :-

$$J = TJ = T_\mu J \quad \text{and.}$$

$$J' = TJ' = T_{\mu'} J'$$

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

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

By Proposition-1, μ and μ' are proper policies and $J = J_\mu$.
and $J' = J_{\mu'}$

$$\text{Now, } J = TJ = T^2 J = \dots = T^K J.$$

$$\text{Also, } J = TJ \leq T_{\mu'} J \leq T_{\mu'}^2 J \leq \dots \leq T_{\mu'}^K J.$$

$$(J \leq T_\mu, J \leq T_{\mu'}^2 J \leq \dots)$$

Thus, $J = T^K J \leq T_{\mu'}^K J \quad \forall K \geq 1$.

Thus as $K \rightarrow \infty$, we obtain $J \leq \lim_{K \rightarrow \infty} T_{\mu'}^K J = T_{\mu'} J = J'$

A symmetric argument then shows that $J' \leq J$. In other words $J = J'$ or $J_\mu = J_{\mu'}$.

Hence "T" has almost one fixed point.

* We now show that T has at least one fixed point:-

let μ be a proper policy. By Assumption (A), \exists a proper policy.

Suppose μ' be another policy s.t.

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

$$\begin{aligned} (T_{\mu'} J_\mu)(i) &= \sum_{j \in S} p_{ij} (\mu'(i)) \\ &\quad (g(i, \mu'(i), j) + J_{\mu'}(j)) \\ &= \min_{\mu \in A(i)} \sum_{j \in S} p_{ij} (\mu) (g(i, \mu, j) + J_\mu(j)) \\ &= (T J_\mu)(i). \end{aligned}$$

Then, $T_\mu J_\mu \geq T J_\mu = T_{\mu'} J_\mu$

Since $J_\mu \geq T_{\mu'} J_\mu$, by Proposition-1, μ' is a proper policy and $J_\mu \geq T_{\mu'} J_\mu \geq T_{\mu''} J_\mu \geq \lim_{K \rightarrow \infty} T_{\mu^K} J_\mu = J_{\mu''}$.

Repeating this argument, we obtain a seq $\{\mu_k\}$ s.t. each μ_k is proper and $J_{\mu_k} = T_{\mu_k} J_\mu \geq T J_{\mu_k} = T_{\mu_{k+1}} J_{\mu_k} \geq \lim_{n \rightarrow \infty} T_{\mu_n} J_\mu = J_{\mu_{k+1}}$.



Since the set of proper policies is finite, some policy μ will get repeated in the sequence.

$$\text{Thus, } J_\mu = TJ_\mu$$

$\Rightarrow J_\mu$ is a fixed point of T and by the uniqueness property, it is the unique fixed pt. of T .

* Next we show, that $J_\mu = J^*$ and $T^k J \rightarrow J^*$ as $k \rightarrow \infty$
 $\forall J \in \mathbb{R}^n$.

Let $e = (1, 1, \dots, 1)$ & $\delta > 0$ be a scalar.

Also let $\hat{J} \in \mathbb{R}^n$ be a s.t $T_\mu \hat{J} = \hat{J} - \delta e$

rewrite

(Note :- $\hat{J} = T_\mu \hat{J} + \delta e$)

$$\hat{J} = \underbrace{g_\mu + \delta e}_{\downarrow} + P_\mu \hat{J}$$

Single state cost.

f is the cost vector corresponding to μ with g_μ replaced with $g_\mu + \delta e$.

Since μ is proper, \hat{J} is unique. further $J_\mu \leq \hat{J}$.

$$\text{Thus, } J_\mu = TJ_\mu \leq T\hat{J} \leq T_\mu \hat{J} = \hat{J} - \delta e \leq \hat{J}$$

$$\Rightarrow J_\mu = T^k J_\mu \leq T^k \hat{J} \leq T^{k-1} \hat{J} \leq \dots \leq \hat{J} \quad \forall k \geq 1.$$

Thus, $\{T^k \hat{J}\}$ forms a bounded monotone sequence and it converges,
 $T^k f \rightarrow \hat{J}$ as $k \rightarrow \infty$ for some $\hat{J} \in \mathbb{R}^n$

$$T\tilde{J} = T(\lim_{K \rightarrow \infty} T^K J) = \lim_{K \rightarrow \infty} T^{K+1} J = \tilde{J}$$

We know that J_μ is a unique fixed point of T , i.e.,
 $J_\mu = TJ_\mu$.

Since \tilde{J} is also a fixed pt., it must happen that
 $\tilde{J} = J_\mu$.

Also, $J_\mu - \delta e = TJ_\mu - \delta e \leq \underbrace{T(J_\mu - \delta e)}_{\text{seen in a prev. class.}} \leq TJ_\mu = J_\mu$.

$$\Rightarrow T(J_\mu - \delta e) \leq T^2(J_\mu - \delta e).$$

Thus $T^K(J_\mu - \delta e)$ is monotonically increasing and bounded above.

Also, $\lim_{K \rightarrow \infty} T^K(J_\mu - \delta e) = J_\mu$.

For any J , we can find $\delta > 0$ s.t.

$$J_\mu - \delta e \leq J \leq \tilde{J}$$

cost-to-go vector for policy μ but with single stage costs $\gamma \mu + \delta e$.

By monotonicity of T ,

$$T^K(J_\mu - \delta e) \leq T^K J \leq T^K \tilde{J}, \quad K \geq 1$$

In the limit,

$$J_\mu = \lim_{K \rightarrow \infty} T^K(J_\mu - \delta e) \leq \lim_{K \rightarrow \infty} T^K J \leq \lim_{K \rightarrow \infty} T^K \tilde{J} \leq \tilde{J}_\mu.$$

To show that $J_\mu = J^*$, take any policy $\pi = \{\mu_0, \mu_1, \dots\}$

Then, $T_{\mu_0} T_\mu, \dots, T_{\mu_k}, J_0 \geq T^* J_0$.

Hence J_0 is an arbitrary vector. Taking \limsup as $k \rightarrow \infty$ on either side, we obtain

$J^* \geq J_\mu \Rightarrow J_\mu$ is optimal.

$\Rightarrow \mu$ is optimal stationary policy and $J_\mu = J^*$

(c) Note: If μ is optimal, then $J_\mu = J^*$.

By Assumption (A)-(B), μ is proper.

Thus by prop. 1, $T_\mu J^* = T_\mu J_\mu = J_\mu = J^* = TJ^*$.

Conversely let $J^* = TJ^* =$ Then μ is proper and $J^* = J_\mu$

$\Rightarrow \mu$ is optimal.

9/2/2023

Lec-10

* Recall:-

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

$i = 1, \dots, n.$

we show that T is a contraction w.r.t. some norm.

\exists a vector $\xi = (\xi(1) - \xi(n))$. s.t. $\xi(i) > 0$.

$\forall i = 1 \dots n$ and a scalar $\beta \in (0, 1)$ s.t.

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

and $\|TJ - T\bar{J}\|_\xi \leq \beta \|J - \bar{J}\|_\xi$

$\forall J, \bar{J} \in \mathbb{R}^n$.

Aside:-

$$(TJ) = ((TJ)(1), \dots, (TJ)(n))$$

$$(T\bar{J}) = ((T\bar{J})(1), \dots, (T\bar{J})(n))$$

we say that T is a contraction w.r.t. Some norm

$$\|\cdot\|_\xi \text{ if } \|TJ - T\bar{J}\|_\xi \leq \alpha \|J - \bar{J}\|_\xi$$

for some $\alpha \in (0, 1)$.

* Proposition-3:- Assume all stationary policies are proper.
Then \exists a vector ξ with $\xi(i) > 0 \forall i$
s.t. the mappings T and T_μ are contractions w.r.t. $\|\cdot\|_\xi$.

Proof:-

Consider a new SSPP where transition probabilities are same as original but transition costs are all equal to -1 except that upon termination costs are 0 .

Let $\hat{J}(i)$ = optimal cost to go from state i in the new problem.

Since, $J^* = TJ^*$, thus.

$$\begin{aligned}\hat{J}(i) &= -1 + \min_{\mu \in A(i)} \sum_{j=1}^n p_{ij}(\mu) \cdot \hat{J}(j), \\ &\leq -1 + \sum_{j=1}^n p_{ij}(\mu(i)) \cdot \hat{J}(j).\end{aligned}$$

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

Then $\forall i$, $\xi_p(i) \geq 1$ and for all stationary μ .

$$-\hat{J}(i) \geq 1 + \sum_{j=1}^n p_{ij}(\mu(i)) (-\hat{J}(j)).$$

$$\text{or } \xi_p(i) \geq 1 + \sum_{j=1}^n p_{ij}(\mu(i)) \xi_p(j).$$

$$\text{or } \sum_{j=1}^n p_{ij}(\mu(i)) \xi_p(j) \leq \xi_p(i) - 1 \leq \beta \cdot \xi_p(i), \forall i.$$

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

Now for any stationary policy μ , state i , vectors $J, \bar{J} \in \mathbb{R}^n$.

$$\begin{aligned}|(T_\mu J)(i) - (T_\mu \bar{J})(i)| &= \left| \sum_{j=1}^n p_{ij}(\mu(i)) (J(j) - \bar{J}(j)) \right| \\ &\leq \sum_{j=1}^n p_{ij}(\mu(i)) |J(j) - \bar{J}(j)| \\ &\leq \left(\sum_{j=1}^n p_{ij}(\mu(i)) \xi_p(j) \right) \underbrace{\left(\max_{j=1, \dots, n} \frac{|J(j) - \bar{J}(j)|}{\xi_p(j)} \right)}_{||J - \bar{J}||_{\xi_p}} \\ &\leq \beta \xi_p(i) ||J - \bar{J}||_{\xi_p}\end{aligned}$$

$$\Rightarrow \frac{|(T_\mu J)(i) - (T_{\mu}\bar{J})(i)|}{\varepsilon_p(i)} \leq \beta \|J - \bar{J}\|_{\varepsilon_p}$$

$\forall i = 1, \dots, n.$

$$\therefore \max_{i=1, \dots, n} \left(\frac{|(T_\mu J)(i) - (T_{\mu}\bar{J})(i)|}{\varepsilon_p(i)} \right) \leq \beta \|J - \bar{J}\|_{\varepsilon_p}$$

or

$$\|T_\mu J - T_{\mu}\bar{J}\|_{\varepsilon_p} \leq \beta \|J - \bar{J}\|_{\varepsilon_p}$$

Note also that

$$(T_\mu J)(i) \leq (T_{\mu}\bar{J})(i) + \beta \varepsilon_p(i) \max_{j=1, \dots, n} \left(\frac{|J(j) - \bar{J}(j)|}{\varepsilon_p(j)} \right)$$

Taking min over μ on both sides

$$(TJ)(i) \leq (T\bar{J})(i) + \beta \varepsilon_p(i) \max_{j=1, \dots, n} \left(\frac{|J(j) - \bar{J}(j)|}{\varepsilon_p(j)} \right)$$

Interchanging J and \bar{J} , we also obtain :-

$$(T_{\mu}\bar{J})(i) \leq (T_\mu J)(i) + \beta \varepsilon_p(i) \max_{i=1, \dots, n} \left(\frac{|J(i) - \bar{J}(i)|}{\varepsilon_p(i)} \right)$$

Again taking min over μ on either side :-

$$(T\bar{J})(i) \leq (TJ)(i) + \beta \varepsilon_p(i) \max_{j=1, \dots, n} \left(\frac{|J(j) - \bar{J}(j)|}{\varepsilon_p(j)} \right)$$

It is also easy to see that:-

$$\|TJ - T\bar{J}\|_2 \leq \beta \|J - \bar{J}\|_2,$$

* Numerical approaches for solving Bellman equation

1. Value Iteration :-

(a.) Start with an arbitrary $J_0 = J \in \mathbb{R}^n$,

(b.) Iterate as $J_{k+1} = TJ_k$, $k=0, 1, 2, \dots$

In other words,

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

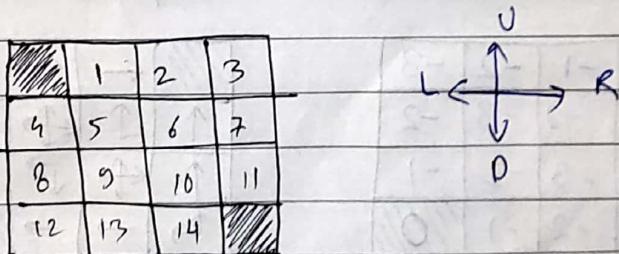
Now, $J_k \rightarrow J^*$ as $k \rightarrow \infty$

where

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

$i=1, \dots, n$.

Ex:- Grid World (Sutton's Ch. 4).



* Non-terminal state = {1, 2, ..., 13}.

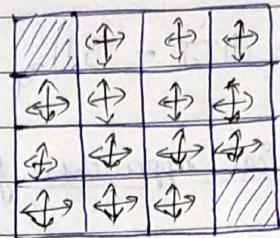
* Terminal state = [shaded square]

* Rewards = -1 is non-terminal states.

Apply value iteration :-

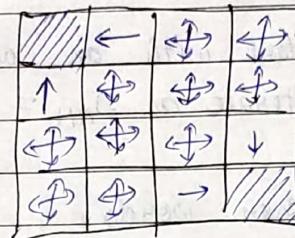
$R=0$

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0



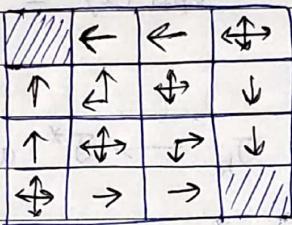
$k=1$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0



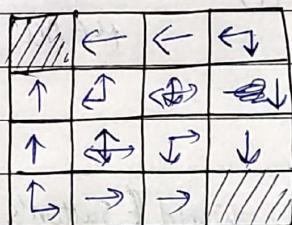
$R=2$

0	-1	-2	-2
-1	-2	-2	-2
-2	-2	-2	-1
-2	-2	-1	0



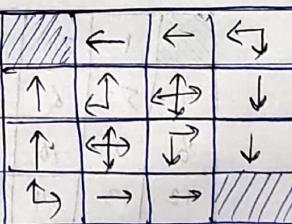
$k=3$

0	-1	-2	-3
-1	-2	-3	-2
-2	-3	-2	-1
-3	-2	-1	0



$R=4$

0	-1	-2	-3
-1	-2	-3	-2
-2	-3	-2	-1
-3	-2	-1	0



Converged.

* Policy iteration :- Iterate over policies.

updates policies via a 2-step procedure

outer loop

{ policy update (policy improvement).

Inner loop

{ find value of new policy (policy evaluation)

}

Go to outer loop.

$T \circ T = T$

?

* Procedure :-

1. Start with a proper policy μ_0 .

2. Policy evaluation :- Given policy μ_k , Compute $J^{\mu_k}(i)$,
 $i \in S$.

as the solⁿ to:

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

$i = 1, \dots, n$

3. Policy improvement : Find a new policy μ_{k+1} s.t.

$$\mu_{k+1}(i) = \underset{\mu \in A(i)}{\operatorname{argmin}} \left\{ \sum_{j=1}^n p_{ij}(u) [g(i, u, j) + J^{\mu_k}(j)] \right\}$$

$i = 1, \dots, n$

$$(\text{OR}) T_{\mu_{k+1}} J^{\mu_k} = T J^{\mu_k}$$

$$(T_{\mu_{k+1}} J^{\mu_k})^{(i)} = \sum_{j=1}^n p_{ij}(\mu_{k+1}(i)) [g(i, \mu_{k+1}(i), j) + J^{\mu_k}(j)]$$

$$= \min_{u \in A(i)} \sum p_{ij}(u) [g(i, u, j) + J^{\mu_k}(j)]$$

Lec:

Result:

Suppose $T: \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a contraction.

Thus, $\|Tx - Ty\| \leq \alpha \|x - y\|$, for any $x, y \in \mathbb{R}^d$ and $0 < \alpha < 1$.

Note that,

$$\|T^2x - T^2y\| \leq \alpha \|Tx - Ty\| \leq \alpha^2 \|x - y\|$$

where $T^2 = T \circ T$

Proceeding similarly, $\|T^n x - T^n y\| \leq \alpha^n \|x - y\|$

Thus the sequence $\{T^n x\}$ is cauchy. Thus it has a limit \bar{y}

$$\Rightarrow T^n x \rightarrow \bar{y} \text{ as } n \rightarrow \infty$$

We say that $\{x_n\}$ is cauchy if given any $\epsilon > 0$

$$\exists N > 0, \forall m, n > N \quad \|x_m - x_n\| < \epsilon$$

$$\|x_m - x_n\| < \epsilon$$

$$\Rightarrow \lim_{n \rightarrow \infty} T^n x = \bar{y}$$

applying T on both sides we obtain :-

$$T\left(\lim_{n \rightarrow \infty} T^n x\right) = T\bar{y}$$

$$\lim_{n \rightarrow \infty} T^{n+1}x = T\bar{y} \quad (\text{by continuity of } T)$$

$\underbrace{\hspace{1cm}}_{\bar{y}}$

$\Rightarrow T\bar{y} = \bar{y}$ or that \bar{y} is a fixed pt. of T .

Suppose \bar{z} is one more fixed pt. of T , i.e.

$$T\bar{z} = \bar{z}. \text{ By contraction property}$$

$$\|\bar{z} - T\bar{y}\| \leq \alpha \|\bar{z} - \bar{y}\|$$

Since, $T\bar{z} = \bar{z}$ & $T\bar{y} = \bar{y}$, we have

$$\|\bar{z} - \bar{y}\| \leq \alpha \|\bar{z} - \bar{y}\|$$

Cannot happen unless $\bar{z} = \bar{y}$

$\Rightarrow T$ has a unique fixed pt.

Another way,

$$\|T^n x - T^n y\| \leq \alpha^n \|x - y\|$$

Suppose we have \bar{y} in place of $y \Rightarrow T^n \bar{y} = \bar{y}$

$$\Rightarrow \|T^n x - \bar{y}\| \leq \alpha^n \|x - \bar{y}\| \Rightarrow T^n x \rightarrow \bar{y} \text{ as } n \rightarrow \infty$$

Policy Iteration

Procedure :-

1. Start with a proper policy μ_0 .

2. Policy evaluation :- Given policy μ_k , complete compute $J^{\mu_k}(i)$, $\forall i \in S$, as the solution to

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

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

3. Policy improvement :- Find a new policy μ_{k+1} s.t.

$$\mu_{k+1}(i) = \underset{\mu \in A(i)}{\operatorname{argmin}} \sum_{j \in S} p_{ij}(\mu) (g(i, \mu, j) + J^{\mu_k}(j)), \quad i \in S.$$

In other words, we find μ_{k+1} s.t.

$$T_{\mu_{k+1}} J^{\mu_k} = T J^{\mu_k}.$$

Go back to (2), find $J^{\mu_{k+1}}$ and repeat process if $J^{\mu_{k+1}}(i) < J^{\mu_k}(i)$ for at least one $i \in S$. If $J^{\mu_{k+1}}(i) = J^{\mu_k}(i)$, $\forall i \in S$, then stop and output the optimal policy-value pair as (μ_k, J^{μ_k}) .

* Proposition :- The policy iteration algorithm generates an improving sequence of proper policies, i.e., $J^{K+1}(i) \leq J^K(i) \forall i \in S$, and ∇K . Further, it terminates with an optimal policy in a finite no. of steps.

Proof :- Given a proper policy μ , new policy $\bar{\mu}$ is optimal as $T_{\bar{\mu}}J^\mu = TJ^\mu$. Then,

$$J^\mu = T_\mu J^\mu \geq TJ^\mu = T_{\bar{\mu}}J^\mu$$

$$\Rightarrow J^\mu \geq T_{\bar{\mu}}J^\mu$$

Repeatedly applying $T_{\bar{\mu}}$ on either side we get,

$$J^\mu \geq T_{\bar{\mu}}J^\mu \geq T_{\bar{\mu}}^2J^\mu \geq \dots \geq T_{\bar{\mu}}^KJ^\mu \geq \dots \geq \lim_{K \rightarrow \infty} (T_{\bar{\mu}}^KJ^\mu)$$

$$\Rightarrow J^{\bar{\mu}}(i) \leq J^\mu(i) \quad \forall i \in S.$$

Since μ is proper, it follows that $\bar{\mu}$ is also proper.
since $\exists i \in S$ s.t. $J^{\bar{\mu}}(i) = \infty$.

$\Rightarrow J^{\bar{\mu}}(i) = \infty$ which cannot happen since μ is proper.

If μ is not optimal, then, $J^{\bar{\mu}}(i) < J^{\mu(i)}$ for at least one $i \in S$.

$$\text{Else, } J^\mu = J^{\bar{\mu}} = T_{\bar{\mu}}J^{\bar{\mu}} = T_{\bar{\mu}}J^\mu = TJ^\mu = T_{\bar{\mu}}J^{\bar{\mu}}$$

$$\Rightarrow J^\mu = TJ^\mu$$

$$\& J^{\bar{\mu}} = TJ^{\bar{\mu}}$$

$$\Rightarrow J^\mu = J^{\bar{\mu}} = J^*$$

In this case, μ is optimal and so is $\bar{\mu}$ & the optimal cost is $J^\mu = J^{\bar{\mu}}$.

Since, no. of proper policies is finite, convergence will happen in a finite no. of steps.

★ Modified PT :-

- Start with μ_0 .
- PE :- Run value iteration for a-period chosen, m_k steps at iteration k .
i.e. Solve $J_k = T_{\mu_k}^{m_k} J_k$ for given m_k .
- PI :- Find μ_{k+1} s.t.

$$T_{\mu_{k+1}} J_k = T J_k.$$

★ Syllabus for Mid-term I :-

- (15 Marks) - Intro to RL, Multi-armed Bandits (Sutton +
(4 ques). Barto),
(3+1). - Finite horizon problems (Bertsekas, Vol. I)
Optimal Control & DP.) (Ch.1)
- SSPP (Bertsekas, Vol. II \rightarrow Ch.2)