

1.1. Minimization of total cost - Introduction.

the stationary discrete-time dynamic system:

$$* x_{k+1} = f(x_k, u_k, w_k) \quad k=0, 1, 2, \dots$$

$x_k \in S$, $u_k \in C$, $w_k \in D$, D is a countable set.

$$* u_k \in U(x_k) \subseteq C, w_k \sim P(\cdot | x_k, u_k) \text{ ind of } k.$$

* $P(w_k | x_k, u_k)$ is the probability w_k occur, w_k ind of w_0, \dots, w_{k-1}

* Given x_0 , we want to find a policy $\Pi = \{\mu_0, \mu_1, \dots\}$

$$\mu_k: S \rightarrow C, \mu_k(x_k) \in U(x_k), \forall x_k \in S, k=0, 1, 2, \dots$$

that minimize the cost function.

$$J_\Pi(x_0) = \lim_{N \rightarrow \infty} E_{w_k} \left[\sum_{k=0}^{N-1} \alpha^k g(x_k, \mu_k(x_k), w_k) \right].$$

$\alpha \in (0, 1)$: discount factors.

* Π : the set of all admissible policies π .

* Optimal cost function: $J^*(x) = \min_{\pi \in \Pi} J_\Pi(x). \quad \forall x \in S$

* An optimal policy for a given state x is one that attains the optimal cost $J^*(x)$. This policy depends on x ,

but sometimes ind. of x .

* Stationary: π is stationary if $\pi = \{\mu_0, \mu_1, \dots\}$ which is referred to as the stationary policy.

* μ is optimal if $J_\mu(x) = J^*(x)$.

1.1.1. The Finite-Horizon DP Algorithm.

Consider any admissible policy $\Pi = \{\mu_0, \mu_1, \dots\}$, any positive integer, any function $J: S \rightarrow \mathbb{R}$.

Suppose we accumulate the costs of the first N stage, and we add to them some terminal cost of the form $\alpha^N J(x_N)$, where J is some function for a total expected cost.

$$E_{w_k} \left[\alpha^N J(x_N) + \sum_{k=0}^{N-1} \alpha^k g(x_k, \mu_k(x_k), w_k) \right].$$

$$k=0, \dots, N-1$$

The minimum cost can be calculated by starting with $\alpha^N J(x)$ and by carrying out N iterations of DP algorithm:

$$J_{N-k}(x) = \min_{u \in U(x)} E \{ \alpha^{N-k} g(x, u, w) + J_{N-k+1}(f(x, u, w)) \}.$$

with initial condition: $J_N(x) = \alpha^N J(x)$.

$J_{N-k}(x)$: optimal cost of the last k stages starting from state x .

$J_0(x)$: the optimal N -stage cost.

More conveniently. we denote $v_k(x) \triangleq \frac{J_{N-k}(x)}{\alpha^{N-k}}$.

then $v_N(x) = J_0(x)$. rewrite DP Algorithm as.

$$v_{k+1}(x) = \min_{u \in U(x)} E \{ g(x, u, w) + \alpha v_k(f(x, u, w)) \}, \quad k=0, 1, \dots, N-1.$$

$$v_0(x) = J(x).$$

1.1.2. Shorthand Notation and Monotonicity.

$\forall J: S \rightarrow \mathbb{R}$. consider DP mapping to J

$$(TJ)(x) = \min_{u \in U(x)} E \{ g(x, u, w) + \alpha J(f(x, u, w)) \}.$$

TJ : optimal cost function for the one-stage problem that has stage cost g and terminal cost αJ

$\forall J: S \rightarrow \mathbb{R}$. control $\mu: S \rightarrow C$

$$(T_\mu J)(x) = E \{ g(x, \mu(x), w) + \alpha J(f(x, \mu(x), w)) \}.$$

cost function associated with μ .

T^k : composition of the mapping T with itself k times.

$$(T^k J)(x) = (T(T^{k-1} J))(x)$$

$$(T^0 J)(x) = J(x), \quad x \in S.$$

$$\text{Similarly: } (T_\mu^k J)(x) = (T_\mu(T_\mu^{k-1} J))(x)$$

$$(T_\mu^0 J)(x) = J(x).$$

$(T^K J)$: optimal cost for the k -stage, α -discounted problem with initial state x , cost per stage g and terminal cost function $\alpha^K J$.

Similarly $(T_\mu^K J)$: the cost of a policy $\{\mu_0, \mu_1, \dots\}$ for the same problem.

Consider k -stage policy $\pi = \{\mu_0, \dots, \mu_{k-1}\}$.

$(T_{\mu_0} T_{\mu_1} \dots T_{\mu_{k-1}} J)(x)$ is defined recursively.

$$(T_{\mu_0} T_{\mu_1} \dots T_{\mu_{k-1}} J)(x) = (T_{\mu_0} (T_{\mu_1} \dots T_{\mu_{k-1}} J))(x)$$

represents the cost under π with initial state x , cost per stage g , terminal cost function $\alpha^K J$.

*Lemma 1.1.1 (Monotonicity Lemma).

For any functions $J: S \rightarrow \mathbb{R}$ and $J': S \rightarrow \mathbb{R}$

$$J(x) \leq J'(x) \quad \forall x \in S.$$

for any stationary policy $\mu: S \rightarrow C$. we have

$$(T^K J)(x) \leq (T^K J')(x). \quad \forall x \in S, k=1, 2, \dots$$

$$(T_\mu^K J)(x) \leq (T_\mu^K J')(x). \quad \forall x \in S, k=1, 2, \dots$$

Define unit function $e: S \rightarrow \mathbb{R}$. $e(x)=1$. $\forall x \in S$.

$$\forall \text{ scalar } r, (T(J+r e))(x) = (TJ)(x) + \alpha r$$

$$(T_\mu(J+r e))(x) = (T_\mu J)(x) + \alpha r$$

* Lemma 1.1.2. $\forall k, J: S \rightarrow \mathbb{R}$, stationary policy μ and scalar r ,

$$(T^K(J+r e))(x) = (T^K J)(x) + \alpha^K r$$

$$(T_\mu^K(J+r e))(x) = (T_\mu^K J)(x) + \alpha^K r$$

1.1.3. A Preview of Infinite Horizon Results.

a) convergence of the DP Algorithm.

J_0 denote zero function: $J_0(x) = 0 \quad \forall x \in S$.

Since the infinite horizon cost of a policy is the limit of its k -stage cost as $k \rightarrow \infty$. Hopefully,

$$J^*(x) = \lim_{k \rightarrow \infty} (T^k J_0)(x) \quad x \in S.$$

Also, $\alpha < 1$ and J bounded. $\alpha^k J$ diminishes with k , then if $\alpha < 1$. $J^*(x) = \lim_{k \rightarrow \infty} (T^k J)(x)$. value iteration starts with any J .

b) Bellman's Equation.

Since by definition, we have for all $x \in S$.

$$(T^{k+1} J_0)(x) = \min_{u \in U(x)} E_w \{ g(x, u, w) + \alpha (T^k J_0)(f(x, u, w)) \}.$$

It is reasonable to speculate that if $\lim_{k \rightarrow \infty} T^k J_0 = J$, take limits on both sides of above equation:

$$J^*(x) = \min_{u \in U(x)} E_w \{ g(x, u, w) + \alpha J^*(f(x, u, w)) \}.$$

$\Rightarrow J^* = TJ^*$: J^* is a fixed point of T .

\nearrow
Bellman's Equation

c) characterization of Optimal Stationary Policy.

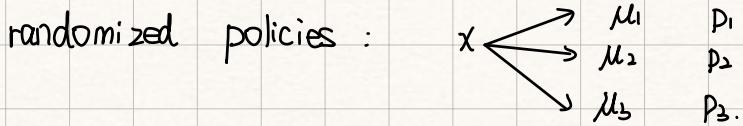
view Bellman's Equation $J^* = TJ^*$ as DP Algorithm
if $\mu(x)$ attains minimum on TJ^* for all x , then the stationary policy is optimal.

1.1.4. Randomized and history dependent policies.

At each time k . μ is applied to z_k . μ is Markov because

μ is independent of x_0, \dots, x_{k-1} .

denote $h_k = \{x_0, \mu_0, x_1, \mu_1, \dots, x_k\}$. Prob



* Proposition 1.1.1 (Adequacy of Markov Policies).

Assume Control space is countable. Initial state distribution takes value over a countable set

The probability distribution of each pair (x_k, u_k) and the expected cost of each stage corresponding to a randomized-history-dependent policy can also be obtained with a randomized Markov policy.

1.2. Discounted problems with bounded cost per stage.

Assumption D. (not very restrictive)

The cost per stage g satisfies $|g(x, u, w)| \leq M$, $\forall (x, u, w) \in S \times C \times D$. M is scalar. $0 < \alpha < 1$

If. S, C, D are finite, then g is always bounded.

DP Algorithm converges to the optimal cost function for an arbitrary bounded starting function J .

$$\lim_{K \rightarrow \infty} E \left\{ \sum_{k=0}^K \alpha^k g(x_k, \mu_k(x_k), w_k) \right\} \rightarrow 0 \quad \text{as } N \rightarrow \infty.$$

* Proposition 1.2.1: (Convergence of DP Algorithm).

For any bounded function $J: S \rightarrow \mathbb{R}$, the optimal cost function satisfies:

$$J^*(x) = \lim_{N \rightarrow \infty} (T^N J)(x). \quad \forall x \in S.$$

proof. \forall positive integer k , $x_0 \in S$ (initial states). policy

$\pi = \{\mu_0, \mu_1, \dots\}$. we break down the cost $J^*(x_0)$ into

the portions incurred over the first k stage and over the remaining stage.

$$\begin{aligned} J_\pi(x_0) &= \lim_{N \rightarrow \infty} E \left\{ \sum_{x=0}^{N-1} \alpha^k g(x_k, \mu_k(x_k), w_k) \right\} \\ &= \lim_{N \rightarrow \infty} E \left\{ \sum_{k=0}^{K-1} \alpha^k g(x_k, \mu_k(x_k), w_k) \right\} + \lim_{N \rightarrow \infty} E \left\{ \sum_{x=K}^{N-1} \alpha^k \right. \\ &\quad \left. g(x_k, \mu_k(x_k), w_k) \right\}. \end{aligned}$$

By Assumption D, $|g(x_k, \mu_k(x_k), w_k)| \leq M$.

$$\left| \lim_{N \rightarrow \infty} E \left\{ \sum_{x=K}^{N-1} \alpha^k g(x_k, \mu_k(x_k), w_k) \right\} \right| \leq M \sum_{k=K}^{\infty} \alpha^k = \frac{\alpha^K M}{1-\alpha}$$

$$\begin{aligned} J_\pi(x_0) - \frac{\alpha^K M}{1-\alpha} - \alpha^K \max_{x \in S} |J(x)| &\leq E \left\{ \alpha^K J(x_K) + \sum_{k=0}^{K-1} \alpha^k g(x_k, \mu_k(x_k), w_k) \right\} \\ &\leq J_\pi(x_0) + \frac{\alpha^K M}{1-\alpha} + \alpha^K \max_{x \in S} |J(x)| \end{aligned}$$

Take minimum over π . we obtain for all x_0, k .

$$\begin{aligned} J^*(x_0) - \frac{\alpha^K M}{1-\alpha} - \alpha^K \max_{x \in S} |J(x)| &\leq (T^K J)(x_0) \\ &\leq J^*(x_0) + \frac{\alpha^K M}{1-\alpha} + \alpha^K \max_{x \in S} |J(x)|. \end{aligned}$$

let $k \rightarrow \infty \Rightarrow J^*(x_0) = (T^K J)(x_0)$.

Gives us a way to compute J^* approximately.

For a given stationary policy μ . we can compute its approximation

* Corollary 1.1.1. For every stationary policy μ , the associated cost function satisfies

$$J_\mu(x) = \lim_{N \rightarrow \infty} (T_\mu^N J)(x), \quad \forall x \in S.$$

J^* is the unique solution to the Bellman's equation:

* Proposition 1.2.2. The optimal cost function J^* satisfies

$$J^*(x) = \min_{u \in U(x)} E_w \{ g(x, u, w) + \alpha J^*(f(x, u, w)) \}. \text{ or } J^* = T J^*$$

Furthermore, J^* is the unique solution of this equation within the class of bounded functions.

proof. $\forall x \in S, N$

$$J^*(x) - \frac{\alpha^N M}{1-\alpha} \leq (T^N J_0)(x) \leq J^*(x) + \frac{\alpha^N M}{1-\alpha}$$

J_0 is the zero function,

apply T to this relation and using the Monotonicity Lemma

1.1.1 and 1.1.2, we obtain $\forall x \in S, \forall n$,

$$(TJ^*)(x) - \frac{\alpha^{N+1} M}{1-\alpha} \leq (T^{N+1} J_0)(x) \leq (TJ^*)(x) + \frac{\alpha^{N+1} M}{1-\alpha}$$

Take $N \rightarrow \infty$, and $\lim_{N \rightarrow \infty} (T^{N+1} J_0)(x) = J^*(x)$

$$\Rightarrow (TJ^*)(x) \leq J^*(x) \leq (TJ^*)(x) \Rightarrow J^* = TJ^*$$

To show uniqueness, if J is bounded and $J = TJ$, then

$$J = \lim_{N \rightarrow \infty} T^N J, \text{ by Prop 1.2.1 we have } J = J^*.$$

* Corollary 1.2.2.1. For every stationary policy μ .

$$J_\mu(x) = E_w \{ g(x, \mu(x), w) + \alpha J_\mu(f(x, \mu(x), w)) \}$$

$$\text{or } J_\mu(x) = T_\mu J_\mu.$$

Furthermore, J_μ is the unique solution of this equation within the class of bounded functions.

* Proposition 1.2.3. (Necessary and Sufficient Condition for Optimality).

A stationary policy μ is optimal if and only if $\mu(x)$ attains the minimum in Bellman's Equation ($J^* = TJ^*$) for each $x \in S$ i.e. $TJ^* = T_\mu J^*$

proof: \Rightarrow If $TJ^* = T_\mu J^*$. By Bellman's Equation

$J^* = TJ^* = T_\mu J^*$, which is unique solution of Bellman's equation : $J^* = J_\mu$. and μ is optimal

\Leftarrow If stationary policy μ is optimal, $J^* = J_\mu$. By corollary 1.2.2.1,

$$J^* = T_\mu J^*$$

Combining this with Bellman's Equation $J^* = TJ^*$

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

* Proposition 1.2.4. For any two bounded functions

$$J: S \rightarrow \mathbb{R} . \quad J': S \rightarrow \mathbb{R} \quad \forall k=0,1,\dots$$

$$\max_{x \in S} |(T^K J)(x) - (T^K J')(x)| \leq \alpha^K \max_{x \in S} |J(x) - J'(x)|$$

proof. Denote $C = \max_{x \in S} |J(x) - J'(x)|$

$$\Rightarrow J(x) - C \leq J'(x) \leq J(x) + C$$

Apply T^K in this relation and use Monotonicity lemma 1.1.1, 1.1.2

$$(T^K J)(x) - \alpha^K C \leq (T^K J')(x) \leq (T^K J)(x) + \alpha^K C \quad x \in S$$

$$\Rightarrow \max_{x \in S} |(T^K J)(x) - (T^K J')(x)| \leq \alpha^K C \quad x \in S$$

* Corollary 1.2.4.1. $J, J': S \rightarrow \mathbb{R}$. \forall stationary policy μ .

$$\max_{x \in S} |(T^K_\mu J)(x) - (T^K_\mu J')(x)| \leq \alpha^K \max_{x \in S} |J(x) - J'(x)|, \quad K=0,1,\dots$$

1.3. Finite-State System — Computational Methods

the state, control and disturbance spaces are finite sets

DP \Rightarrow control of a finite-state Markov chain

Notations:

State space $S = \{1, 2, \dots, n\}$

transition probability : $P_{ij}(u) = P(x_{k+1}=j | x_k=i, u_k=u) \quad i,j \in S, u \in U(i)$

update function : $x_{k+1} = f(x_k, u_k, w_k)$

disturbance probability : $P(\cdot | x, u)$.

$$P_{ij}(u) = P(w_{ij}(u) | i, u)$$

where $w_{ij}(u) \triangleq \{w \in D | f(i, u, w)=j\}$.

To simplify, assume cost per stage does not depend on w

$$(g(i, u, w) = g(i, u))$$

$$\text{define } g(i, u) = \sum_{j=1}^n P_{ij}(u) \tilde{g}(i, u, j)$$

$\tilde{g}(i, u, j)$: the cost of starting from i , use control u , transite to j .

$g(i, u)$: expected cost at state i , use control u .

$$TJ(i) = \min_{u \in U(i)} [g(i, u) + \alpha \sum_{j=1}^n P_{ij}(u) J(j)]$$

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

$TJ, T_\mu J$ can be represented by the n -dimensional vectors.

$$J = \begin{bmatrix} J(1) \\ \vdots \\ J(n) \end{bmatrix} \quad TJ = \begin{bmatrix} (TJ)(1) \\ \vdots \\ (TJ)(n) \end{bmatrix} \quad T_\mu J = \begin{bmatrix} (T_\mu J)(1) \\ \vdots \\ (T_\mu J)(n) \end{bmatrix}$$

For a stationary policy μ , we denote P_μ as transition probability matrix

$$P_\mu = \begin{bmatrix} P_{11}(\mu(1)) & \cdots & P_{1n}(\mu(1)) \\ \vdots & \ddots & \vdots \\ P_{n1}(\mu(n)) & \cdots & P_{nn}(\mu(n)) \end{bmatrix}$$

$$g_\mu: \text{the cost vector } g_\mu = \begin{bmatrix} g(1, \mu(1)) \\ \vdots \\ g(n, \mu(n)) \end{bmatrix}$$

$$\Rightarrow J_\mu = T_\mu J = g_\mu + \alpha P_\mu J_\mu$$

A stationary policy μ . J_μ can be computed

$$(I - \alpha P_\mu) J_\mu = g_\mu \Rightarrow J_\mu = (I - \alpha P_\mu)^{-1} g_\mu$$

$(I - \alpha P_\mu)^{-1}$ exists: $\alpha < 1$. P_μ 's eigenvalues $< 1 \Rightarrow \alpha P_\mu$'s eigenvalues < 1
 $\Rightarrow I - \alpha P_\mu$ is invertible

1.3.1 Value Iteration and Error Bounds.

A n -dimensional vector J , compute J, TJ, \dots successively. By prop 1.2.1.

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

Furthermore, by Prop 1.2.4. error sequence $|(T^K J)(i) - J^*(i)|$ is bounded by

$\alpha^k C$. C is a constant.

This method is called. value iteration.

$$J_\mu(i) = g(i, \mu(i)) + \sum_{k=1}^{\infty} \alpha^k E \{ g(x_k, \mu(x_k)) \mid x_0=i \}.$$

$$\Rightarrow g_\mu + \left(\frac{\alpha \beta}{1-\alpha} \right) e \leq J_\mu \leq g_\mu + \left(\frac{\alpha \bar{\beta}}{1-\alpha} \right) e.$$

where $e = (1, 1, \dots, 1)^T$, $\beta = \min_i g(i, \mu(i))$, $\bar{\beta} = \max_i g(i, \mu(i))$

$$\Rightarrow \left(\frac{\beta}{1-\alpha}\right)e \leq J\mu + \left(\frac{\alpha\beta}{1-\alpha}\right)e \leq TJ \leq J\mu + \left(\frac{\alpha\bar{\beta}}{1-\alpha}\right)e \leq \left(\frac{\bar{\beta}}{1-\alpha}\right)e$$

suppose we have a vector J , compute

$$T\mu J = J\mu + \alpha P_\mu J \quad (1)$$

subtracting this equation from the relation

$$J\mu = J\mu + \alpha P_\mu J\mu \quad (2)$$

$$(1), (2) \Rightarrow J\mu - J = T\mu J - J + \alpha P_\mu (J\mu - J).$$

variational form of the equation $J\mu = T\mu J\mu$.

replace $J\mu$ with $J\mu - J$, $J\mu$ with $T\mu J - J$, we have

$$\begin{aligned} \left(\frac{\gamma}{1-\alpha}\right)e &\leq T\mu J - J + \left(\frac{\alpha\gamma}{1-\alpha}\right)e \\ &< J\mu - J \leq T\mu J - J + \left(\frac{\alpha\bar{\gamma}}{1-\alpha}\right)e \leq \left(\frac{\bar{\gamma}}{1-\alpha}\right)e \end{aligned}$$

where $\gamma = \min_i [(T\mu J)(i) - J(i)]$, $\bar{\gamma} = \max_i [(T\mu J)(i) - J(i)]$

$$\Leftrightarrow J + \frac{\underline{c}}{1-\alpha}e \leq T\mu J + \underline{c}e \leq J\mu \leq T\mu J + \bar{c}e \leq J + \frac{\bar{c}}{1-\alpha}e$$

$$\text{where } \underline{c} = \frac{\alpha\gamma}{1-\alpha}, \bar{c} = \frac{\alpha\bar{\gamma}}{1-\alpha}$$

* Proposition 1.3.1. $\forall J$, state i and k . we have.

$$(T^K J)(i) + \underline{c}_k \leq (T^{K+1} J)(i) + \underline{c}_{k+1}$$

$$\leq J^*(i)$$

$$\leq (T^{K+1} J)(i) + \bar{c}_{k+1} \leq (T^K J)(i) + \bar{c}_k$$

$$\text{where } \underline{c}_k = \frac{\alpha}{1-\alpha} \min_{i=1,\dots,n} [(T^K J)(i) - (T^{K-1} J)(i)]$$

$$\bar{c}_k = \frac{\alpha}{1-\alpha} \max_{i=1,\dots,n} [(T^K J)(i) - (T^{K-1} J)(i)]$$

* Termination Issue - Optimality of the obtained Policy.

$\forall J$, if we compute TJ and policy μ attaining the minimum in the calculation of TJ ($T\mu J = TJ$), then we can obtain the bound on the suboptimality of μ .

$$\max_i [J\mu(i) - J^*(i)] \leq \frac{\alpha}{1-\alpha} (\max_i [(TJ)(i) - J(i)] - \min_i [(TJ)(i) - J(i)]) \quad (3)$$

$$\text{By Prop 1.3.1, } k=1 \quad \underline{c}_1 \leq J^*(i) - (TJ)(i) \leq \bar{c}_1 \quad (4)$$

replacing T with $T\mu$:

$$\underline{c}_1 \leq J\mu(i) - (T\mu J)(i) = J\mu(i) - TJ(i) \leq \bar{c}_1 \quad (5).$$

subtracting the above two equations: (4) - (5) \Rightarrow (3)

In practice, we terminate as $\bar{C}_k - C_k$ is sufficiently small.

One can take the final estimate of J^* the median

$$\tilde{J}_k = T^k J + \left(\frac{\bar{C}_k - C_k}{2} \right) e$$

or the average:

$$\hat{J}_k = T^k J + \frac{\alpha}{n(1-\alpha)} \sum_{i=1}^n (T^k J)(i) - (T^{k-1} J)(i) e$$

Because for eq(3), $\exists \bar{\varepsilon} > 0$, s.t. if μ satisfy

$$\max_i [J_\mu(i) - J^*(i)] < \bar{\varepsilon},$$

then μ is optimal.

let K be such that for all $k > K$, we have

$$\frac{\alpha}{1-\alpha} (\max_i [(T^K J)(i) - (T^{K-1} J)(i)] - \min_i [(T^K J)(i) - (T^{K-1} J)(i)]) < \bar{\varepsilon}$$

then eq(3) implies that $\forall k \geq K$, the stationary policy that attains the minimum in the k th value iteration is optimal.

1.3.3. Policy Iteration.

The policy iteration generates a sequence of stationary policies, each with improved cost over the preceding one. Given stationary policy μ and cost function J_μ , an improved $\{\bar{\mu}, \bar{\mu}, \dots\}$ is computed by minimization in the DP equation corresponding to $J_\mu \leq \bar{T}_{\bar{\mu}} J_\mu = T J_{\bar{\mu}}$, and the process is repeated.

* Proposition 1.3.4. Let μ and $\bar{\mu}$ be stationary policies such that

$$\bar{T}_{\bar{\mu}} J_\mu = T J_\mu, \text{ or equivalently}$$

$$g(i, \bar{\mu}(i)) + \alpha \sum_{j=1}^P p_{ij} (\bar{\mu}(i)) J_\mu(j) = \min_{u \in U(i)} [g(i, u) + \alpha \sum_{j=1}^P p_{ij}(u) J_\mu(j)]$$

Then we have.

$$J_{\bar{\mu}}(i) \leq J_\mu(i)$$

Furthermore, if μ is not optimal, strict inequality holds in the above

equation for at least one state i

proof $J_\mu = T_\mu J_\mu$ and $T_{\bar{\mu}} J_\mu = T J_\mu$, $\forall i$;

$$\begin{aligned} J_{\mu(i)} &= g(i, \mu(i)) + \alpha \sum_{j=1}^n P_{ij}(\mu(i)) J_\mu(j) \\ &\geq g(i, \bar{\mu}(i)) + \alpha \sum_{j=1}^n P_{ij}(\bar{\mu}(i)) J_\mu(j) \\ &= (T_{\bar{\mu}} J_\mu)(i) \end{aligned}$$

Apply $T_{\bar{\mu}}$ on both sides of this inequality and use the monotonicity of $T_{\bar{\mu}}$ and Corollary 1.2.1.1. we have.

$$J_\mu \geq T_{\bar{\mu}} J_\mu \geq \dots \geq T_{\bar{\mu}}^k J_\mu \geq \dots \geq \lim_{N \rightarrow \infty} T_{\bar{\mu}}^N J_\mu = J_{\bar{\mu}}$$

If $J_\mu = J_{\bar{\mu}}$, then $J_\mu = T_{\bar{\mu}} J_\mu$.

By hypothesis $\Rightarrow T_{\bar{\mu}} J_\mu = T J_\mu \Rightarrow J_\mu = T J_\mu$.

* Policy Iteration Algorithm

Step 1 (Initialization) Guess an initial stationary policy μ^0

Step 2. (Policy Evaluation) Given the stationary policy μ^k , compute the corresponding cost function J_{μ^k} from the linear system of equations : $(I - \alpha P_{\mu^k}) J_{\mu^k} = g_{\mu^k}$

Step 3. (Policy Improvement) Obtain a new stationary policy μ^{k+1} satisfying

$$J_{\mu^{k+1}} J_{\mu^k} = T J_{\mu^k}$$

If $J_{\mu^k} = T J_{\mu^k}$, stop. else, return to step 2 and repeat.

