

VALUE ITERATION, POLICY ITERATION: PRELIMINARY RESULTS

VI-PI. 1

Recall

$$V^*(x) = \min_{\mu} \left[\sum_{t=0}^{+\infty} \gamma^t \mathbb{E} \left[g(x_t, \mu(x_t), w_t) \right] \right] \quad x_0 = x$$

optimal policy is stationary

$$V^M(x) = \sum_{t=0}^{+\infty} \gamma^t \mathbb{E} \left[g(x_t, \mu(x_t), w_t) \right]$$

$$\min_{\mu} \left[\mathbb{E} \left[g(x, \mu, w) + \gamma V^*(f(x, \mu, w)) \right] \right]$$

we know that (Bellman's equation)

$$V^*(x) = \min_{\mu} \left[\mathbb{E} \left[g(x, \mu, w_0) \right] + \gamma \sum_{i=1}^M P(i|x, \mu) V^*(i) \right] \quad \forall x \in X$$

$$\Rightarrow T^*[V] : \mathbb{R}^M \rightarrow \mathbb{R}^M \quad \min_{\mu} \left[\mathbb{E} \left[g(x, \mu, w) + \gamma V(f(x, \mu, w)) \right] \right]$$

$$T^*[V]_x = \min_{\mu} \left[\mathbb{E} \left[g(x, \mu, w_0) \right] + \gamma \sum_{i=1}^M P(i|x, \mu) V(i) \right] \quad \forall x \in X$$

$$V^* = \text{solution to } V = T^*[V] \quad (\text{fixed point})$$

similarly, $\forall x \in X$

$$\begin{aligned} V^M(x) &= \sum_{t=0}^{+\infty} \gamma^t \mathbb{E} \left[g(x_t, \mu(x_t), w_t) \right] = \mathbb{E} \left[g(x, \mu(x), w_0) \right] + \gamma \sum_{t=1}^{+\infty} \gamma^{t-1} \mathbb{E} \left[g(x_t, \mu(x_t), w_t) \right] \\ &= \mathbb{E} \left[g(x, \mu(x), w_0) \right] + \gamma V^M(x_{\perp}) = \mathbb{E} \left[g(x, \mu(x), w) \right] + \gamma \sum_{i=1}^M P(i|x, \mu(x)) V^M(i) \end{aligned}$$

$$\Rightarrow T^M[V] : \mathbb{R}^M \rightarrow \mathbb{R}^M \quad T^M[V]_x = \mathbb{E} \left[g(x, \mu(x), w_0) \right] + \gamma \sum_{i=1}^M P(i|x, \mu(x)) V(i) \quad x \in X$$

$$V^M = \text{solution to } V = T^M[V] \quad (\text{fixed point}) \quad \mathbb{E} \left[g(x, \mu(x), w) + \gamma V(f(x, \mu(x), w)) \right]$$

result $V = T^*[V]$ and $V = T^M[V]$... solution exist and is

unique \Leftarrow Banach-Caccioppoli theorem (the tool)

(S, d) = complete metric space $F: S \rightarrow S$ contraction mapping

$$\text{i.e. } \exists \rho < 1 : d(F(s_1), F(s_2)) \leq \rho \cdot d(s_1, s_2) \quad \forall s_1, s_2 \in S$$

$s = F(s)$ always admits one and only one solution

clearly this implies continuity

Proof: uniqueness

By contradiction, suppose that $s_1 = F(s_1)$ and $s_2 = F(s_2)$. Then

$$d(s_1, s_2) = d(F(s_1), F(s_2)) \leq p d(s_1, s_2) < d(s_1, s_2) \dots \text{contradiction!}$$

$p < 1$

existence

let s_0 be an arbitrary point in S and define

$$s_{k+1} = F(s_k) = F^{k+1}(s_0) \quad k = 0, 1, 2, \dots$$

$$\text{Then } d(s_{k+1}, s_k) = d(F(s_k), F(s_{k-1})) \leq p d(s_k, s_{k-1}) \leq \dots \leq p^k d(s_1, s_0)$$

For any $M > m$
triangle ineq.

$$\begin{aligned} d(s_M, s_m) &\leq d(s_M, s_{M-1}) + d(s_{M-1}, s_{M-2}) + \dots + d(s_{m+1}, s_m) \leq \dots \\ &\dots \leq \sum_{k=m}^{M-1} d(s_{k+1}, s_k) \leq \sum_{k=m}^{M-1} p^k d(s_1, s_0) \leq \sum_{k=m}^{\infty} p^k d(s_1, s_0) = \frac{p^m}{1-p} d(s_1, s_0) \end{aligned}$$

so for any $M > m$, $d(s_M, s_m) \xrightarrow{M \rightarrow \infty} 0$, i.e. $\{s_m\}$ is Cauchy convergent

$$\text{completeness} \Rightarrow \exists \bar{s} : s_m \xrightarrow{m \rightarrow \infty} \bar{s}$$

$$\text{However, } \bar{s} = \lim_{k \rightarrow \infty} s_{k+1} = \lim_{k \rightarrow \infty} F(s_k) = F(\bar{s}) \quad \text{i.e. } \bar{s} \text{ is a fixed point.}$$

QED

$s_0 =$ arbitrary

$$s_{k+1} = F(s_k)$$

\Rightarrow fixed-point iteration
(subsequent approximations)

\Rightarrow algorithm to find fixed point
 $\bar{s} = F(\bar{s})$

$$d(s_{k+1}, \bar{s}) = d(F(s_k), F(\bar{s})) \leq p d(s_k, \bar{s}) \leq \dots \leq p^{k+1} d(s_0, \bar{s}) \quad (\text{exponential convergence})$$

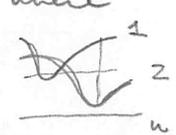
Fundamental property

Both T^* and T^M are contractions from \mathbb{R}^M to \mathbb{R}^M with respect to the max-norm distance

T^* : let V_1 and V_2 two value functions (i.e. two vectors in \mathbb{R}^M)

$$\begin{aligned} \max_{x \in X} \left| T^*[V_1] - T^*[V_2] \right| &= \max_{x \in X} \left| \min_u \left[E[g(x, u, w)] + \gamma \sum_{i=1}^M p(i|x, u) V_1(i) \right] + \right. \\ &\quad \left. - \min_u \left[E[g(x, u, w)] + \gamma \sum_{i=1}^M p(i|x, u) V_2(i) \right] \right| \leq \end{aligned}$$

\leq assume $\min_u [\dots V_2(i)] \leq \min_u [\dots V_1(i)]$ and attained for u_2^*
 (if $\min_u [\dots V_2(i)] \leq \min_u [\dots V_1(i)]$ evaluate everything in u_1^* where $\min_u [\dots V_1(i)]$ is attained)
 \rightarrow distance increases \leftarrow min



$$\leq \max_{x \in X} \left| \mathbb{E}[g(x, u_2^*, w)] + \gamma \sum_{i=1}^M p(i|x, u_2^*) V_1(i) - \mathbb{E}[g(x, u_2^*, w)] - \gamma \sum_{i=1}^M p(i|x, u_2^*) V_2(i) \right|$$

$$= \max_{x \in X} \left| \gamma \sum_{i=1}^M p(i|x, u_2^*) (V_1(i) - V_2(i)) \right| \leq \max_{x \in X} \gamma \sum_{i=1}^M p(i|x, u_2^*) |V_1(i) - V_2(i)|$$

$$\leq \max_{i \in X} |V_1(i) - V_2(i)| \cdot \max_{x \in X} \gamma \sum_{i=1}^M p(i|x, u_2^*) =$$

$$= \gamma \cdot \max_{x \in X} |V_1(x) - V_2(x)|$$

T^μ : same, simpler V_1 and $V_2 \in \mathbb{R}^m$ value functions (ie)

$$\max_{x \in X} |T^\mu[V_1] - T^\mu[V_2]| =$$

$$= \max_{x \in X} \left| \mathbb{E}[g(x, \mu(x), w)] + \gamma \sum_{i=1}^M p(i|x, \mu(x)) V_1(i) - \mathbb{E}[g(x, \mu(x), w)] - \gamma \sum_{i=1}^M p(i|x, \mu(x)) V_2(i) \right|$$

$$= \max_{x \in X} \left| \gamma \sum_{i=1}^M p(i|x, \mu(x)) (V_1(i) - V_2(i)) \right| \leq$$

$$\leq \gamma \max_{x \in X} \sum_{i=1}^M p(i|x, \mu(x)) |V_1(i) - V_2(i)| \leq \gamma \max_{i \in X} |V_1(i) - V_2(i)| \cdot \max_{x \in X} \sum_{i=1}^M p(i|x, \mu(x)) =$$

$$= \gamma \max_{x \in X} |V_1(x) - V_2(x)|$$

Thus, $V = T^*[V]$ and $V = T^\mu[V]$ admits one and only one solution $\rightarrow V^*(x)$ and $V^\mu(x)$

optimal policy \iff optimal value function $V^*(x)$ } key problem: compute $V^*(x), V^\mu(x)$
 policy iteration \iff value function for μ $V^\mu(x)$

VALUE ITERATION AND POLICY ITERATION ALGORITHMS

• $V^*(x) = \min_u \left[\mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^m P(i|x, u) V^*(i) \right] \quad x=1, 2, \dots, m \quad VI-P1.4$

piecewise linear in V^* \Rightarrow system of non-linear eqs

• $V^\mu(x) = \mathbb{E}[g(x, \mu(x), w)] + \gamma \sum_{i=1}^m P(i|x, \mu(x)) V^\mu(i)$

\Downarrow linear in $V^\mu \Rightarrow$ system of linear eqs

tough
(impossible, curse of dimensionality)

$(I - \gamma \cdot P) V^\mu = \left[\mathbb{E}[g(x, \mu(x), w)] \right]_{x=1}^m$ better but still tough if m is large

\Rightarrow use fixed-point iterations! \Rightarrow VALUE ITERATION algorithm (VI)

$V^0(x) \quad x \in X$ arbitrarily chosen (need not be a value function associated to some policy)

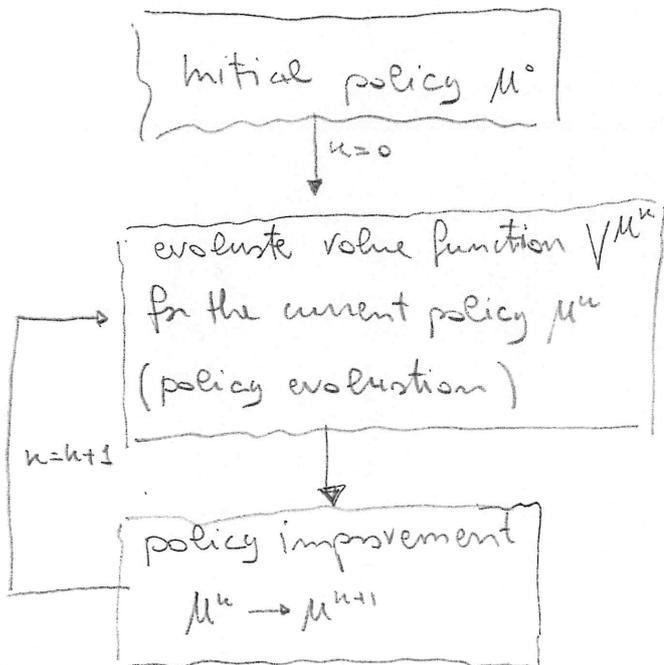
• $V^{k+1}(x) = T^*[V^k(x)]$ converges to $V^*(x)$

• $V^{k+1}(x) = T^\mu[V^k(x)]$ converges to $V^\mu(x)$

exponentially fast

OBS: $V^k(x) = [T^*]^k [V^0(x)] \quad \text{or} \quad V^k(x) = [T^\mu]^k [V^0(x)]$

Alternative to VI to compute $V^* \Rightarrow$ POLICY ITERATION (PI)



• sequence of policies, improving at each iteration ($V^{\mu^{k+1}} \leq V^{\mu^k} \forall x \in X$)

• convergence in a finite no. of iterations can be very large

• faster convergence than VI for large state space

big improvement in the first few iterations

property: MONOTONICITY

VI-PI.4.bis

Suppose $V_1(x) \leq V_2(x) \quad \forall x \in X$. Then, for every μ we have

$$T^\mu[V_1]_x \leq T^\mu[V_2]_x \quad \forall x \in X \quad \text{and} \quad T^*[V_1]_x \leq T^*[V_2]_x \quad \forall x \in X$$

Proof:

$$\begin{aligned} T^\mu[V_1]_x &= \mathbb{E}[g(x, \mu(x), w)] + \gamma \sum_{i=1}^M p(i|x, \mu(x)) V_1(i) \\ &\leq \mathbb{E}[g(x, \mu(x), w)] + \gamma \sum_{i=1}^M p(i|x, \mu(x)) V_2(i) = T^\mu[V_2]_x \end{aligned}$$

$$\mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^M \overset{\geq 0}{p(i|x, u)} V_1(i) \leq \mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^M \overset{\geq 0}{p(i|x, u)} V_2(i)$$

$\forall x \in X \quad \forall u \in U$. Hence,

$$\begin{aligned} T^*[V_1]_x &= \min_{u \in U} \left[\mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^M p(i|x, u) V_1(i) \right] \\ &\leq \min_{u \in U} \left[\mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^M p(i|x, u) V_2(i) \right] = T^*[V_2]_x \end{aligned}$$

Start with an initial policy $\mu^0: X \rightarrow U$ arbitrary

VI-Pl.5

Policy evolution: solve $V^{\mu^k}(x) = \mathbb{E}[g(x, \mu^k(x), w)] + \gamma \sum_{i=1}^M P(i|x, \mu^k(x)) \cdot V^{\mu^k}(i)$
 $\forall x \in X$

\hookrightarrow either solve system of linear eqs. or use VI for V^{μ^k} exact till convergence $V = T^{\mu^k}[V]$

Policy improvement: $\mu^{k+1}(x) = \underset{u \in U}{\operatorname{argmin}} \mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^M P(i|x, u) V^{\mu^k}(i)$

Properties:

1. $V^{\mu^{k+1}}(x) \leq V^{\mu^k}(x) \quad \forall x \in X$ $V^0 = T^{\mu^{k+1}}[V^{\mu^k}]$

Define $\forall x \in X, V^0(x) = \mathbb{E}[g(x, \mu^{k+1}(x), w)] + \gamma \sum_{i=1}^M P(i|x, \mu^{k+1}(x)) \cdot V^{\mu^k}(i)$
 $= \min_{u \in U} \mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^M P(i|x, u) V^{\mu^k}(i)$
 $\leq \mathbb{E}[g(x, \mu^k(x), w)] + \gamma \sum_{i=1}^M P(i|x, \mu^k(x)) V^{\mu^k}(i) = V^{\mu^k}(x)$

so $V^0(x) \leq V^{\mu^k}(x) \quad \forall x \in X$

Then

$V^{h+1} = T^{\mu^{k+1}}[V^h] \quad h=0, 1, \dots \implies V^h \xrightarrow{h \rightarrow \infty} V^{\mu^{k+1}}$

$V^1 = T^{\mu^{k+1}}[V^0] \leq (\text{monotonicity, } V^0 \leq V^{\mu^k}) \leq T^{\mu^{k+1}}[V^{\mu^k}] = V^0$

$V^2 = T^{\mu^{k+1}}[V^1] \leq (\text{monotonicity, } V^1 \leq V^0) T^{\mu^{k+1}}[V^0] = V^1 \leq V^0$

\vdots

$V^h \leq V^0 \leq V^{\mu^k} \implies V^{\mu^{k+1}} = \lim_{h \rightarrow \infty} V^h \leq V^{\mu^k}$

2. no. of ^{stationary} policies is finite ($\leq M^M$) $\implies V^{\mu^{k+1}}(x) < V^{\mu^k}(x)$ for some x

occurs for a finite no. of iterations, then for some k

$V^{\mu^{k+1}}(x) = V^{\mu^k}(x) \quad \forall x \in X \implies V^{\mu^{k+1}} = V^h \quad \forall h$ in previous iterations

$\implies V^{\mu^k}(x) = V^{\mu^{k+1}}(x) = V^0(x) = \min_u \left[\mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^M P(i|x, u) V^{\mu^k}(i) \right]$

i.e. $V^{\mu^n}(x)$ satisfies the Bellman eq. for optimal value $V = P1.6$
 function $\Rightarrow \mu^n = \mu^*$ optimal policy, $V^{\mu^n} = V^*$ optimal
 value function

PI converges in finite time to the optimal policy!

Q-FACTORS (action-value functions)

$$Q^*: X \times U \rightarrow \mathbb{R} \quad Q^\mu: X \times U \rightarrow \mathbb{R}$$

$$Q^*(x, u) = \mathbb{E}[g(x, u, w_0)] + \gamma \sum_{i=1}^M P(i|x, u) V^*(i) \quad x \in X \quad u \in U$$

$$Q^\mu(x, u) = \mathbb{E}[g(x, u, w_0)] + \gamma \sum_{i=1}^M P(i|x, u) V^\mu(i) \quad x \in X \quad u \in U$$

interpretation: total discounted cost achieved when at $t=0$ (initial time)

$x_0 = x$ and $u_t = u$ (action u is used) and then:

: optimal policy μ^ is applied for $t \geq 1$

μ : policy μ is applied for $t \geq 1$

Clearly, thanks to Bellman equations

$$\bullet \quad V^*(x) = \min_{u \in U} Q^*(x, u) \quad \mu^*(x) = \underset{u \in U}{\operatorname{argmin}} Q^*(x, u)$$

$$\bullet \quad V^\mu(x) = Q^\mu(x, \mu(x)) \quad Q^* \text{ summarizes all the info we need to determine the optimal policy}$$

Bellman equations for Q-factors (from definitions and relations above)

$$Q^*(x, u) = \mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^M P(i|x, u) \cdot \min_{v \in U} Q^*(i, v)$$

$$\text{solution to } Q = T^*[Q] \text{ where } T^*[Q] = \mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^M P(i|x, u) \min_v Q(i, v)$$

$$Q^\mu(x, u) = \mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^M P(i|x, u) Q^\mu(i, \mu(i))$$

$$\text{solution to } Q = T^\mu[Q] \text{ where } T^\mu[Q] = \mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^M P(i|x, u) Q(i, \mu(i))$$

property: T^* and T^M (for Q-factors) are contractions with $\forall l=1, 2$ bits
 respect to the $\max_{x,u}$ distance

$Q_1(x,u), Q_2(x,u)$ given.

$$\begin{aligned} \max_{x,u} |T^*[Q_1] - T^*[Q_2]| &= \max_{x,u} \left| \mathbb{E}[\cancel{g(x,u,w)}] + \gamma \sum_{i=1}^M P(i|x,u) \min_v Q_1(i,v) - \mathbb{E}[\cancel{g(x,u,w)}] + \right. \\ &\quad \left. - \gamma \sum_{i=1}^M P(i|x,u) \min_v Q_2(i,v) \right| \leq \max_{x,u} \left| \gamma \sum_{i=1}^M P(i|x,u) \left| \min_v Q_1(i,v) - \min_v Q_2(i,v) \right| \right| \\ &\leq \max_{x,u} \left| \gamma \sum_{i=1}^M P(i|x,u) |Q_1(i, v(i)) - Q_2(i, v(i))| \right| \quad v(i) = \begin{cases} \operatorname{argmin}_v Q_1(i,v) & \text{if } \min_u Q_1(i,u) \leq \min_u Q_2(i,u) \\ \operatorname{argmin}_v Q_2(i,v) & \text{else} \end{cases} \\ &\leq \max_{i,v} |Q_1(i,v) - Q_2(i,v)| \cdot \max_{x,u} \gamma \sum_{i=1}^M P(i|x,u) \\ &= \gamma \cdot \max_{x,u} |Q_1(x,u) - Q_2(x,u)| \end{aligned}$$

$$\begin{aligned} \max_{x,u} |T^M[Q_1] - T^M[Q_2]| &= \max_{x,u} \left| \mathbb{E}[\cancel{g(x,u,w)}] + \gamma \sum_{i=1}^M P(i|x,u) Q_1(i, \mu(i)) - \mathbb{E}[\cancel{g(x,u,w)}] + \gamma \sum_{i=1}^M P(i|x,u) Q_2(i, \mu(i)) \right| \\ &\leq \max_{x,u} \gamma \cdot \sum_{i=1}^M P(i|x,u) |Q_1(i, \mu(i)) - Q_2(i, \mu(i))| \leq \dots \leq \gamma \max_{x,u} |Q_1(x,u) - Q_2(x,u)| \end{aligned}$$

as before

property: T^* and T^M (for Q-factors) are monotonic

$Q_1(x,u) \leq Q_2(x,u) \quad \forall x \in X, u \in U$ Then, $\forall x \in X, u \in U$

$$\begin{aligned} T^*[Q_1](x,u) &= \mathbb{E}[g(x,u,w)] + \gamma \sum_{i=1}^M P(i|x,u) \min_v Q_1(i,v) \\ &\leq \mathbb{E}[g(x,u,w)] + \gamma \sum_{i=1}^M P(i|x,u) \min_v Q_2(i,v) = T^*[Q_2](x,u) \end{aligned}$$

$$\begin{aligned} T^M[Q_1](x,u) &= \mathbb{E}[g(x,u,w)] + \gamma \sum_{i=1}^M P(i|x,u) Q_1(i, \mu(i)) \\ &\leq \mathbb{E}[g(x,u,w)] + \gamma \sum_{i=1}^M P(i|x,u) Q_2(i, \mu(i)) = T^M[Q_2](x,u) \end{aligned}$$

FACT: since T^* and T^M are contractions, the Bellman eqs. for Q-factors univocally determine $Q^*(x,u)$ and $Q^M(x,u)$ always (can be proved otherwise) and shows the validity of the following Value Iteration method for Q-factors.

in a model-based environment $\rightarrow V^*$ and Q^* are equivalent

in a model-free (RL): μ^* = argmin

$$\mu^*(x) = \underset{u}{\operatorname{argmin}} \underbrace{Q^*(x,u)} = \underset{u}{\operatorname{argmin}} \mathbb{E}[g(x,u,w)] + \gamma \sum_{i=1}^M P(i|x,u) V^*(i)$$

↑ that's all we need to compute the optimal policy
↑ model needed besides V^*

Thus, $Q^*(x,u) = \mathbb{E}[g(x,u,w) + \gamma \min_v Q^*(P(x,u,w), v)]$

$$Q^*(x,u) = T^*[Q^*](x,u) = \mathbb{E}[g(x,u,w)] + \gamma \min_v Q^*(P(x,u,w), v)$$

↑ expected value outside
↑ min inside \Rightarrow suitable for incremental computation

$$V^*(x) = T^*[V^*](x) = \min_u [\mathbb{E}[g(x,u,w) + \gamma V^*(P(x,u,w))]]$$

↑ outside \rightarrow not suitable for incremental computation
↑ inside

VALUE ITERATION for Q-factors

$Q^0(x,u) : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ arbitrary

$$Q^{k+1}(x,u) = \mathbb{E}[g(x,u,w)] + \gamma \sum_{i=1}^M P(i|x,u) \cdot \min_v Q^k(i,v) \xrightarrow{k \rightarrow \infty} Q^*(x,u)$$

$$= T^*[Q^k](x,u)$$

$$Q^{k+1}(x,u) = \mathbb{E}[g(x,u,w)] + \gamma \sum_{i=1}^M P(i|x,u) Q^k(i, \mu(i)) \xrightarrow{k \rightarrow \infty} Q^\mu(x,u)$$

$$= T^\mu[Q^k](x,u)$$

either solve sys. of linear equations or use VI

POLICY ITERATION for Q-factors

$\mu^0(x) : \mathcal{X} \rightarrow \mathcal{U}$ initial policy (arbitrary)

policy evaluation: solve $Q^{\mu^k}(x,u) = \mathbb{E}[g(x,u,w)] + \gamma \sum_{i=1}^M P(i|x,u) Q^{\mu^k}(i, \mu^k(i))$

$$Q^{\mu^k}(x,u) = T^{\mu^k}[Q^{\mu^k}](x,u) \quad \forall x \in \mathcal{X} \quad u \in \mathcal{U}$$

policy improvement: $\mu^{k+1}(x) = \underset{u}{\operatorname{argmin}} Q^{\mu^k}(x,u)$

OBS: the fact that

VI-PI:

$$Q^{\mu^u}(x, u) = \mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^m P(i|x, u) Q^{\mu^u}(i, \mu^u(i)) \quad (*)$$

always admit one and only one solution can be also seen by noting that

$$Q^{\mu^u}(x, \mu^u(x)) = \mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^m P(i|x, u) Q^{\mu^u}(i, \mu^u(i)) \quad \text{is the}$$

Bellman equation for value functions, which admits one and only one solution. Hence, (*) univocally determines $Q^{\mu^u}(x, \mu^u(x))$ and

$$Q^{\mu^u}(x, \mu^u(x)) = V^{\mu^u}(x) \quad (\text{solution to } *). \quad \text{However, once } Q^{\mu^u}(x, \mu^u(x))$$

are determined, (*) determines $Q^{\mu^u}(x, u)$ for the other values of u since right-hand side of (*) becomes given.

$$\text{At every policy evolution: } Q^{\mu^u}(x, u) : Q^{\mu^u}(x, \mu^u(x)) = V^{\mu^u}(x)$$

$$\text{and therefore } Q^{\mu^u}(x, u) = \mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^m P(i|x, u) V^{\mu^u}(i)$$

Hence, at every policy improvement

$$\mu^{u+1}(x) = \underset{u}{\operatorname{argmin}} \mathbb{E}[g(x, u, w)] + \gamma \sum_{i=1}^m P(i|x, u) V^{\mu^u}(i)$$

same as PI for value iteration

PI for Q-factors and PI for value function give the same policy at every iteration

convergence and other properties inherited