

Convergence of VI

V_i is pointwise increasing

ξ

Bounded

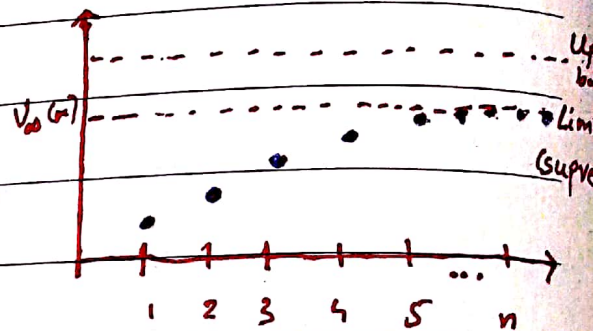
\Downarrow

pointwise convergent

Monotone convergence Theorem

$\forall i, x$

$$V_i(x) \leq V_{i+1}(x)$$



$$V_i(x) \leq V(x) < \infty$$

\downarrow

$$V_\infty(x) = V(x)$$

$\forall x$

Moreover

$$\text{Policy: } \mu_\infty = \mu^*$$



∞ horizon undiscounted optimal control problem

$$\min_{\mu} \mathbb{E} \left[\sum_{k=0}^{\infty} \rho(X_k, U_k) \right], \quad \rho(x, u) \geq 0 \\ \forall x, u$$

$$X_+ \sim P_x(X_+ | x, u)$$

$$U_+ \sim P_{\mu}^{\mu(x)}(u | x)$$

parameters of policy model ' $\mu(x)$ ' depend upon state, x .

$$\mu : X \rightarrow \Pi$$

space of policy
PDF parameters

VI algorithm: $\forall x$

$$\mu_i(x) := \arg \min_{\mu} \left\{ \mathbb{E}_{U \sim P_{\mu}^{\mu(x)}} \left[\rho(x, U) + V_i(X_+^U) \right] \right\}$$

$$V_{i+1} := \mathbb{E}_{U \sim P_{\mu_i}^{\mu_i(x)}} \left[\rho(x, U) + V_i(X_+^U) \right]$$

η arbitrary, "admissible" policy

$$J(x|\eta) = \mathbb{E}_{U_k \sim p_U^\eta(x_k)} \left[\sum_{k=0}^{\infty} \rho(x_k, U_k) \mid x_0 = x \right]$$

should be finite!

→ What properties comprise an admissible policy?

Comparing iteration by iteration

$$V_{i+1} := \mathbb{E}_{U \sim p_U^{\mu(x)}} \left[\rho(x, U) + V_i(x, U) \right]$$

$$\Lambda_{i+1} := \mathbb{E}_{U \sim p_U^{\eta(x)}} \left[\rho(x, U) + \Lambda_i(x, U) \right]$$

$$\Lambda_0 = 0$$

①

$$V_1(x) = \mathbb{E}_{U \sim p_U^{\mu_0(x)}} \left[\rho(x, U) \right]$$

①

$$\Lambda_1(x) = \mathbb{E}_{U \sim p_U^{\eta(x)}} \left[\rho(x, U) \right]$$

where

$$\mu_0(x) = \operatorname{argmin}_{\sigma} \mathbb{E}_{U \sim p_U^{\sigma}} \left[\rho(x, U) \right]$$

Because of minimizer (or optimality)
 $P_0(x)$

$$V_1(x) \leq \Lambda_1(x)$$

(2)

$$V_2(x) = \mathbb{E}_{U \sim P_U^{P_1(x)}} \left[p(x, U) + V_1(x_+^U) \right]$$

$$= \mathbb{E}_{U \sim P_U^{P_1(x)}} \left[p(x, U) + \mathbb{E}_{U_+ \sim P_{U_+}^{P_0(x_+^U)}} [p(x_+^U, U_+)] \right]$$

$$= \mathbb{E}_{\substack{U \sim P_U^{P_1(x)} \\ U_+ \sim P_{U_+}^{P_0(x_+^U)}}} \left[p(x, U) + p(x_+^U, U_+) \right]$$

$$P_1(x) = \underset{Q}{\operatorname{argmin}} \mathbb{E}_{\substack{U \sim P_U^Q \\ U_+ \sim P_{U_+}^{P_0(x_+^U)}}} \left[p(x, U) + p(x_+^U, U_+) \right]$$

Due to optimality of $P_1(x)$

$$\mathbb{E}_{\substack{U \sim P_U^{P_1(x)} \\ U_+ \sim P_{U_+}^{P_0(x_+^U)}}} \left[p(x, U) + p(x_+^U, U_+) \right] \leq \mathbb{E}_{\substack{U \sim P_U^{\eta(x)} \\ U_+ \sim P_{U_+}^{P_0(x_+^U)}}} \left[p(x, U) + p(x_+^U, U_+) \right]$$

can be cho

$$U_+ \sim P_{U_+}^{\eta(x_+^U)}$$

We already established the argument with the optimality of μ_0 i.e.

$$\mathbb{E}_{U \sim P_U^{\mu_0(x)}} [p(x, U)] \leq \mathbb{E}_{U \sim P_U^{\eta(x)}} [p(x, U)]$$

Stepping one step changes nothing

$$\mathbb{E}_{U \sim P_U^{\mu_0(x_+^U)}} [p(x_+^U, U)] \leq \mathbb{E}_{U \sim P_U^{\eta(x_+^U)}} [p(x_+^U, U)]$$

So, now

$$V_2(x) \leq \Lambda_2(x) \quad \forall x$$

Supposing:

$$V_i(x) \leq \Lambda_i(x) \quad \forall x$$

$$\mu_i(x) = \arg \min_{\mu} \left\{ \mathbb{E}_{U \sim P_U^{\mu(x)}} [p(x, U) + V_i(x_+^U)] \right\}$$

$$V_{i+1}(x) := \mathbb{E}_{U \sim P_U^{\mu_i(x)}} [p(x, U) + V_i(x_+^U)]$$

$$\Lambda_{i+1}(x) = \mathbb{E}_{U \sim P_U^{\eta(x)}} [p(x, U) + \Lambda_i(x_+^U)]$$

show: $V_{i+1}(x) \leq \Lambda_{i+1}(x)$

By optimality of π_i

$$\mathbb{E}_{U \sim P_U^{\pi_i(x)}} [r(x, U) + V_i(x_U^U)] \leq \mathbb{E}_{U \sim P_U^{\pi(x)}} [r(x, U) + V_i(x_U^U)]$$

$$\mathbb{E}_{U \sim P_U^{\pi(x)}} [r(x, U) + V_i(x_U^U)] \leq \mathbb{E}_{U \sim P_U^{\pi}} [r(x, U) + \Lambda_i(x_U^U)]$$

$V_{i+1}(x)$

$<$

$\Lambda_{i+1}(x)$

Proof by induction

Using optimal policy instead of arbitrary

② Showing the boundedness $V_i(x) \leq V(x)$ ↖ optimum cost to go

$$V_i(x) \leq \mathbb{E}_{U_{+k} \sim P_U^{\pi^*(x_{+k})}} \left[\sum_{k=0}^{\infty} r(x_{+k}, U_{+k}) \right]$$

$$\mathbb{E}_{U_{+k} \sim P_U^{\pi^*(x_{+k})}} \left[\sum_{k=0}^{i-1} r(x_{+k}, U_{+k}) \mid x = x \right]$$

adding non zero rewards
cannot be bigger

⊙ Pointwise increasing

$$V_i(x) \leq V_{i+1}(x)$$

$$V_{i+1}(x) = \mathbb{E}_{U \sim P_U^{\mu_i(x)}} [r(x, U) + V_i(x_t^U)]$$

$$V_i(x) = \mathbb{E}_{U \sim P_U^{\mu_{i-1}(x)}} [r(x, U) + V_{i-1}(x_t^U)]$$

μ_{i-1} is minimizer so

$$\mathbb{E}_{U \sim P_U^{\mu_{i-1}(x)}} [r(x, U) + V_{i-1}(x_t^U)]$$

$$\leq \mathbb{E}_{U \sim P_U^{\mu_i(x)}} [r(x, U) + V_{i-1}(x_t^U)]$$

or any policy for that matter

$$V_i(x) \leq \mathbb{E}_{U \sim P_U^{\mu_i(x)}} [r(x, U) + V_{i-1}(x_t^U)]$$

$$| X=x$$

$$V_i(x) = \mathbb{E}_{U_{+k} \sim p_U^{M_{i-1-k}(X_{+k})}} \left[\sum_{i=0}^{i-1} p(X_{+k}, U_{+k}) \right]$$

$$V_{i+1}(x) = \mathbb{E}_{U_{+k} \sim p_U^{M_{i+2-k}(X_{+k})}} \left[\sum_{i=0}^i p(X_{+k}, U_{+k}) \mid X=x \right]$$

Verifying

We have seen this :

$$i=2$$

$$V_2(x) = \mathbb{E}_{\substack{U \sim p_U^{M_1(x)} \\ U_+ \sim p_U^{M_0(x_+)}}} \left[p(x, U) + p(x_+, U_+) \right]$$

$\hookrightarrow U_+ \sim p_U^{M_{i-1-k}(X_{+1})} \Rightarrow 2-1-1 = p_U^{M_0(x_+)}$

$$V_3(x) = \mathbb{E}_{\substack{U \sim p_U^{M_2(x)} \\ U_+ \sim p_U^{M_1(x_+)} \\ U_{+2} \sim p_U^{M_0(x_{+2})}}} \left[p(x, U) + p(x_+, U_+) + p(x_{+2}, U_{+2}) \right]$$

\downarrow
 $U_{+k} \sim p_U^{M_{i+2-k}(X_{+k})}$

\downarrow
 $\sum_{k=0}^i p(X_{+k}, U_{+k})$

$$i=3$$

$$V_i(x) = V_{i-1}(x_t^u)$$

$$V_{i-1}(x_t^u) \leq V_i(x_t^u)$$

$$V_i(x) \leq \mathbb{E}_{u \sim p_u^{V_i(x)}} [f(x, u) + V_{i-1}(x_t^u)]$$

$$\leq \mathbb{E}_{u \sim p_u^{V_i(x)}} [f(x, u) + V_i(x_t^u)]$$

$$\leq V_{i+1}(x)$$