

Convergence of tabular methods

Let's consider the following undiscounted ∞ -horizon problem:

$$\max_{\rho} \mathbb{E} \left[\sum_{k=0}^{\infty} r(X_k, U_k) \right]$$

$$\text{s.t. } X_{k+1} \sim f(x_{k+1} | x_k, u_k)$$

$$U_k = \rho(X_k)$$

Let's assume $r \leq 0$

Let's do tabular VI:

$$\forall x \in \mathcal{X} \quad \rho_i(x) := \arg \max_u \{ r(x, u) + \mathbb{E} [V_i(X_+)] \}$$

$$V_{i+1}(x) := r(x, \rho_i(x)) + \mathbb{E} [V_i(X_+^{\rho_i(x)})]$$

Reminder: $X_+^u \sim f(x_+ | x, u)$

Let's consider $\gamma : \mathcal{X} \rightarrow \mathcal{U}$, an arbitrary policy

Its value reads:

$$V^\eta(x) = E_f \left[\sum_{u=0}^{\infty} r(X_u, \eta(X_u)) \mid X_0 = x \right]$$

Since V^η may fail to converge, we consider so-called admissible policies.

Admissible $\eta \Rightarrow V^\eta > -\infty$

Let's consider generic iterations of the kind:

$$\Lambda_{i+1}(x) := r(x, \eta(x)) + E \left[\Lambda_i(X_i^{(x)}) \right]$$

Start with $\Lambda_0 \equiv 0$

Let's benchmark our VI scheme using Λ_i .

Iteration $i = 1$:

$$V_1(x) = r(x, \rho_0(x))$$

let's say ρ_0 was chosen so as to satisfy $\forall x \in \mathcal{X} \quad r(x, \rho_0(x)) \geq r(x, \eta(x))$

At the same time,

$$\Lambda_1(x) = r(x, \eta(x))$$

$$\Rightarrow V_1 \geq \Lambda_1$$

Iteration 2:

$$V_2(x) = r(x, \rho_1(x)) + \mathbb{E} \left[V_1 \left(X_+^{\rho_1(x)} \right) \right]$$

$$= r(x, \rho_1(x)) + \mathbb{E} \left[r \left(X_+^{\rho_2(x)}, \rho_0(X_+^{\rho_2(x)}) \right) \right]$$

$$\rho_1(x) = \arg \max_u \left\{ r(x, u) + \mathbb{E} \left[r \left(X_+^u, \rho_0(X_+^u) \right) \right] \right\}$$

So since ρ_1 is a maximizer,

$$r(x, \rho_1(x)) + \mathbb{E} \left[r \left(X_+^{\rho_0(x)}, \rho_1(X_+^{\rho_0(x)}) \right) \right] \geq$$

$$r(x, \eta(x)) + \mathbb{E} \left[r \left(X_+^{\rho_0(x)}, \eta(X_+^{\rho_0(x)}) \right) \right] \geq$$

$$r(x, \eta(x)) + \mathbb{E} \left[r \left(X_+^{\eta(x)}, \eta(X_+^{\eta(x)}) \right) \right]$$

So, we may conclude that

$$V_2 \geq A_2$$

Let's assume $V_i \geq A_i$

Consider

$$\rho_i(x) = \arg \max_u \{r(x, u) + \mathbb{E}[V_i(X_+^{\rho_i(x)})]\}$$

$$V_{i+1}(x) = r(x, \rho_i(x)) + \mathbb{E}[V_i(X_+^{\rho_i(x)})]$$

$$A_{i+1}(x) = r(x, \eta(x)) + \mathbb{E}[A_i(X_+^{\eta(x)})]$$

By optimality of ρ_i ,

$$r(x, \rho_i(x)) + \mathbb{E}[V_i(X_+^{\rho_i(x)})] \geq$$

$$r(x, \eta(x)) + \mathbb{E}[V_i(X_+^{\eta(x)})]$$

By the induction hypothesis,

$$r(x, \eta(x)) + \mathbb{E}[V_i(X_+^{\eta(x)})] \geq$$

$$r(x, \eta(x)) + \mathbb{E}[A_i(X_+^{\eta(x)})]$$

Which means,

$$V_{i+1}(x) \geq \Lambda_{i+1}(x), \quad \forall x \in \mathcal{X}$$

So, a proof by induction is complete.
Since this holds for any π (admissible), it also holds for the optimal policy π^* !

The benchmark iterations read:

$$\Lambda_{i+1}(x) = r(x, \pi^*(x)) + \mathbb{E}[\Lambda_i(\tilde{X}_+^{\pi^*(x)})]$$

Unwrapping the recursion yields:

$$\Lambda_{i+1}(x) = \mathbb{E}\left[\sum_{k=0}^{i-1} r(X_k, \pi^*(X_k)) \middle| X_0 = x\right]$$

Adding the infinite tail, and keeping in mind that $r \leq 0$,

$$\begin{aligned} \Lambda_i(x) &\geq \mathbb{E}\left[\sum_{k=0}^{\infty} r(X_k, \pi^*(X_k)) \middle| X_0 = x\right] \\ &= V^*(x) \end{aligned}$$

Looking at two successive iterations:

$$V_{i+1}(x) = r(x, \beta_i(x)) + \mathbb{E}\left[V_i\left(X_t^{\beta_i(x)}\right)\right]$$

$$V_i(x) = r(x, \beta_{i-1}(x)) + \mathbb{E}\left[V_{i-1}\left(X_t^{\beta_{i-1}(x)}\right)\right]$$

Remember, β_{i-1} was optimal on its iteration:

$$\beta_{i-1}(x) = \arg \max_u \{r(x, u) + \mathbb{E}[V_{i-1}(X_t^u)]\}$$

We thus have

$$r(x, \beta_{i-1}(x)) + \mathbb{E}\left[V_{i-1}\left(X_t^{\beta_{i-1}(x)}\right)\right] \geq$$

$$r(x, \beta_i(x)) + \mathbb{E}\left[V_{i-1}\left(X_t^{\beta_i(x)}\right)\right]$$

So,

$$V_i(x) \geq r(x, \beta_i(x)) + \mathbb{E}\left[V_{i-1}\left(X_t^{\beta_i(x)}\right)\right]$$

Let's do some more unwrapping:

$$V_i(x) = \mathbb{E} \left[\sum_{k=0}^{i-1} r(X_k, p_{i-1-k}(X_k)) \middle| X_0 = x \right]$$

$$V_{i-1}(x) = \mathbb{E} \left[\sum_{k=0}^{i-2} r(X_k, p_{i-2-k}(X_k)) \middle| X_0 = x \right]$$

Now shift the starting state
to X_+^u :

$$V_i(X_+^u) = \mathbb{E} \left[\sum_{k=1}^i r(X_k, p_{i-k}(X_k)) \middle| X_0 = X_+^u \right]$$

$$V_{i-1}(X_+^u) = \mathbb{E} \left[\sum_{k=1}^{i-1} r(X_k, p_{i-1-k}(X_k)) \middle| X_0 = X_+^u \right]$$

Add the i th reward to the last
expression to get:

$$\mathbb{E} \left[\sum_{k=1}^{i-1} r(X_k, p_{i-1-k}(X_k)) \middle| X_0 = X_+^u \right] \geq$$

$$\mathbb{E} \left[\sum_{k=1}^{i-1} r(X_k, p_{i-1-k}(X_k)) + r(X_i, p_{i-1}(X_i)) \middle| X_0 = X_+^u \right] =$$

$$\mathbb{E} \left[\sum_{k=1}^i r(X_k, p_{i-k}(X_k)) \middle| X_0 = X_+^u \right] =$$

$$V_i(X_+^u)$$

Therefore,

$$\begin{aligned} V_i(x) &\geq r(x, \rho_i(x)) + E[V_{i+1}(X_+^{\rho_i(x)})] \\ &\geq r(x, \rho_i(x)) + E[V_i(X_+^{\rho_i(x)})] \\ &= V_{i+1}(x) \end{aligned}$$

So, we have :

$$\forall i, x \in \mathcal{X} \quad V_i(x) \geq V_{i+1}(x)$$

$$V_i(x) \geq V^*(x) > -\infty$$

To show convergence, apply
the monotone convergence theorem

$$\text{So, } V_i \rightarrow V_\infty$$

$$\text{We know } V_\infty \geq V^*(x)$$

Since V^* is the optimal value function, it has to be the greatest, i.e., $V^*(x) \geq V_\infty(x), \forall x \in \mathcal{X}$,

so $V_\infty = V^*!$