

AVD894 - Reinforcement Learning

Assignment - I

17-09-2025

SAURABH KUMAR
SC22BL46

Proving convergence of the value iteration algorithm for policy evaluation.

$$\text{Defn: } v^\pi(s) = \sum_a \pi(a|s) \sum_h p(h|s,a) h + \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) \cdot v^\pi(s')$$

In the matrix form,

$$\bar{V}^\pi = \bar{\pi}^\pi + \gamma P^\pi \bar{V}^\pi$$

To compute $v^\pi(s)$, we have the value iteration algorithm that starts with an estimate $v_0(s) \forall s \in S$ which is initialized to 0.

Then, for every k we define

$$\bar{v}_{k+1}^\pi = \bar{\pi}^\pi + \gamma P^\pi \bar{v}_k^\pi$$

To prove that :

$$\lim_{k \rightarrow \infty} \bar{v}_k^\pi = v^\pi, \forall s \in S$$

Define the error, $\delta_k = \bar{v}_k - v^\pi$.

Substitute $\bar{v}_{k+1} = \delta_{k+1} + v^\pi$ and $\bar{v}_k = \delta_k + v^\pi$

$$\text{into } \bar{v}_{k+1} = \bar{\pi}^\pi + \gamma P^\pi \bar{v}_k$$

$$\Rightarrow \delta_{k+1} + v^\pi = \bar{\pi}^\pi + \gamma P^\pi (\delta_k + v^\pi)$$

$$\Rightarrow \delta_{k+1} = -v^\pi + \bar{\pi}^\pi + \gamma P^\pi \delta_k + \gamma P^\pi v^\pi$$

$$= \gamma P^\pi \delta_k - v^\pi + (\underbrace{\bar{\pi}^\pi + \gamma P^\pi v^\pi}_{\bar{v}^\pi})$$

$$= \gamma P^\pi \delta_k$$

$$= \gamma^2 P^\pi \delta_{k-1} = \dots = \gamma^{k+1} P^\pi \delta_0.$$

As $0 \leq P_\pi^k \leq 1 \forall k$ and $\gamma < 1$, so $\gamma^k \rightarrow 0$,

hence $\delta_{k+1} = \gamma^{k+1} P^\pi \delta_0 \rightarrow 0$ as $k \rightarrow \infty$.

As $\delta_k \rightarrow 0$, $v_k \rightarrow v^\pi$.

Contraction Mapping:

$$\text{Take } f(v_1) = \mu^\pi + \gamma P^\pi v_1^\pi$$

$$f(v_2) = \mu^\pi + \gamma P^\pi v_2^\pi$$

$$|f(v_1) - f(v_2)| = |\gamma P^\pi v_1^\pi - \gamma P^\pi v_2^\pi|$$

$$\leq \gamma \|v_1^\pi - v_2^\pi\|, \text{ as } P^\pi \text{ sums to 1.}$$

Thus, for $0 \leq \gamma < 1$, the mapping $\bar{v}(.)$ to $\bar{v}_{k+1}(.)$ is a contraction map.

Because the mapping is a contraction, iterating from any starting v_0 gives

$$\begin{aligned} \|v_k - v^\pi\|_\infty &\leq \gamma \|v_{k+1} - v^\pi\|_\infty \\ &\leq \gamma^k \|v_0 - v^\pi\|_\infty \end{aligned}$$

$$\therefore \|v_k - v^\pi\|_\infty \rightarrow 0 \text{ as } k \rightarrow \infty,$$

hence $v_k \rightarrow v^\pi$ as $k \rightarrow \infty$.