

Value Iteration, Policy Iteration and iLQR

1. Derive the formula to get the minimum number of iterations of Value Iteration that are needed if we want an error on the quality of the policy that is at most ϵ

Value Iteration is a method used to compute the Optimal Value V^* of a Markov Decision Process, the main idea behind this method is to use mathematical properties to find a way to obtain V^* starting from an initial value V_0 .

The reason why we need alternative ways to compute the optimal Value is that the direct computation is highly expensive: given the definition of $V^\pi(s) := r + \gamma \mathbb{E}_{p(\cdot|s,\pi)}[V^\pi(s')]$ it's possible to extract the matrix form of the equation $V^\pi = R + \gamma P^\pi V^\pi$ and as a consequence the value of V^π

$$V^\pi = (I - \gamma P^\pi)^{-1} R$$

But this calculation – even if it's correct – is not feasible since $I - \gamma P^\pi \in \mathbb{R}^{S \times S}$, where S is the number of states of the MDP, and just inverting the matrix will be a $\mathcal{O}(S^3)$ and this evaluation should be done for each possible policy π in order to determine V^* .

In Value Iteration, using two mathematical concepts:

- **Fixed Point**, given a map f we say that x is a fixed point for f if $x = f(x)$.
- **Contraction**, a map $f : X \rightarrow X$ is called a contraction if $\exists k \in (0, 1) \mid \forall (x, x') \in X \times X \rightarrow |f(x) - f(x')| \leq k|x - x'|$.

an iterative algorithm is defined to compute V^* . In fact, those concepts taken together guarantee the following:

Given $f : X \rightarrow X$ a contraction mapping, with a fixed point \bar{x} and an initial value $x_0 \in X$.

Then the iterative algorithm $x_i = f(x_{i-1})$ converges: $x_i \xrightarrow{i \rightarrow \infty} \bar{x}$.

So, we can define Value Iteration by taking $f = BV(s) = \max_a(r(s, a) + \gamma \mathbb{E}_{p(\cdot|s,a)}[V(s')])$ – the Bellman Operator – that, as shown in class, is indeed a contraction and its fixed value is V^* .

As for the number of iteration needed to obtain an error of the policy's quality to be less than ϵ , we want $\|V^{\pi_i} - V^*\| \leq \epsilon$, where π_i is computed using the Bellman Operator for Q :

$$\text{Operator } TQ(r, a) = r(s, a) + \gamma \mathbb{E}_{p(\cdot|s,a)} \max_{a'} [Q(s', a')]$$

$$\text{The policy } \pi_i(s) = \arg\max_a Q_i(s, a) \text{ and } Q_i = TQ_{i-1}$$

It can be shown, proof will be at the end of the section, that $\|V^{\pi_i} - V^*\| \leq \frac{2\gamma^i}{1-\gamma} \|Q_0 - Q^*\|$, so we can rewrite the problem as $\frac{2\gamma^i}{1-\gamma} \|Q_0 - Q^*\| \leq \epsilon$.

$$\begin{aligned} \frac{2\gamma^i}{1-\gamma} \|Q_0 - Q^*\| &= \text{Assume initialization } Q_0 = 0 \Rightarrow = \frac{2(1-(1-\gamma))^i}{1-\gamma} \|Q^*\| \leq \\ &\leq \text{since } 1+x \leq e^x \text{ and } Q^* \in [0, \frac{1}{1-\gamma}] \leq \frac{2e^{-(1-\gamma)i}}{(1-\gamma)^2} \end{aligned}$$

So to have that the error on the quality of the policy that is at most ϵ we have to find the value of i such that $\frac{2e^{-(1-\gamma)i}}{(1-\gamma)^2} \leq \epsilon$, since this would imply $\frac{2\gamma^i}{1-\gamma} \|Q_0 - Q^*\| \leq \epsilon$.

$$\begin{aligned} \frac{2e^{-(1-\gamma)i}}{(1-\gamma)^2} \leq \epsilon &\Leftrightarrow e^{-(1-\gamma)i} \leq \frac{\epsilon(1-\gamma)^2}{2} \Leftrightarrow \log(e^{-(1-\gamma)i}) \leq \log\left(\frac{\epsilon(1-\gamma)^2}{2}\right) \Leftrightarrow \\ &\Leftrightarrow -i(1-\gamma) \leq -\log\left(\frac{2}{\epsilon(1-\gamma)^2}\right) \Leftrightarrow i \geq \frac{\log\left(\frac{2}{\epsilon(1-\gamma)^2}\right)}{1-\gamma} \end{aligned}$$

Proof.

- $\|V^{\pi_i} - V^*\| \leq \frac{2\gamma^i}{1-\gamma} \|Q_0 - Q^*\|:$

$$\begin{aligned} \|V^{\pi_i} - V^*\| &= \max_s -(V^{\pi_i}(s) - V^*(s)) \text{ [this quantity is } \geq 0 \text{ since } V^* \geq V^{\pi_i}] \\ V^{\pi_i}(s) - V^*(s) &= Q^{\pi_i}(s, \pi^i(s)) - Q^*(s, \pi^i(s)) + Q^*(s, \pi^i(s)) - Q^*(s, \pi^*(s)) = \\ &= \gamma \mathbf{E}_{s' \sim P(s, \pi^i(s))} (V^{\pi_i}(s') - V^*(s')) + Q^*(s, \pi^i(s)) - Q^*(s, \pi^*(s)) \geq \\ &[\text{Add the negative quantity: } -Q^{\pi_i}(s, \pi^i(s)) + Q^{\pi_i}(s, \pi^*(s))] \\ &\geq \gamma \mathbf{E}_{s' \sim P(s, \pi^i(s))} (V^{\pi_i} - V^*) + Q^*(s, \pi^i(s)) - Q^{\pi_i}(s, \pi^i(s)) + Q^{\pi_i}(s, \pi^*(s)) - Q^*(s, \pi^*(s)) \geq \\ &[\text{Since } |Q^i(s) - Q^*(s)| \leq \|Q^i - Q^*\|_\infty \leq \gamma^i \|Q^0 - Q^*\|_\infty] \\ &\geq \gamma \mathbf{E}_{s' \sim P(s, \pi^i(s))} (V^{\pi_i} - V^*) - 2\gamma^i \|Q^0 - Q^*\|_\infty \geq \dots \geq -\frac{2\gamma^i}{1-\gamma} \|Q^0 - Q^*\|_\infty \end{aligned}$$

2. Compute $V_{k+1}(s_6)$ following the Value Iteration algorithm.

- States $\{s_i : i \in \{1, \dots, 7\}\}$
- Rewards $r(s, a) = \begin{cases} 0.5 & \text{if } s = s_1 \\ 5 & \text{if } s = s_7 \\ 0 & \text{otherwise} \end{cases}$
- Dynamics $p(s_6|s_6, a_1) = 0.3$ $p(s_7|s_6, a_1) = 0.7$
- Policy $\pi(s) = a_1, \forall s$
- $v_k = [0.5, 0, 0, 0, 0, 0, 5]$
- $\gamma = 0.9$

We want to compute one iteration of the iterative algorithm used to compute the Value Function that, as state before, has a matricial form:

$$V_t = R + \gamma P V_{t-1}$$

In this specific case we have to compute $V_{k+1}(s_6)$ so the components needed to find this quantity are $r(s_6, a_1)$, v_k and s_{k+1} .

- $R = [0.5, 0, 0, 0, 0, 0, 5] \Rightarrow r(s_k = s_6, a_1) = 0$
- $V_{t-1} = v_k = [0.5, 0, 0, 0, 0, 0, 5]$
- $s_{k+1} \sim p(\cdot | s_k = s_6, a_1) = [0, 0, 0, 0, 0, 0.3, 0.7]$

So we can easily compute $v_{k+1}(s_6)$ as

$$r(s_k = s_6, a_1) + \gamma < v_k, p(\cdot | s_k = s_6, a_1) >$$

$$\Downarrow$$

$$v_{k+1}(s_6) = 0 + 0.9(0 * 0.3 + 5 * 0.7) = 3.15$$