goal: derive algorithms to solve MDP/SOCP when model is "known"

refs:

# Neuro-Dynamic Programming

Dimitri P. Bertsekas and John N. Tsitsiklis

chapter 2

---

o given finite MDP/SOCP $(X, U, P, c)$ — infinite-horizon, exponentially-discounted

i.e. $\min_{u} E[c(x, u)]$        $c(x, u) = \sum_{t=0}^{\infty} \gamma^t \cdot \mathcal{L}(x_t, u_t)$

s.t. $x^+ \sim P(x, u)$, $|X|, |U| < \infty$        $\gamma \in (0, 1)$

o we can define Bellman operators associated with MDP, $\pi : X \to \Delta(U)$

$T : \mathbb{R}^X \to \mathbb{R}^X$   — nonlinear / piecewise-affine

$: v \longmapsto (Tv)(x) = \min_{u \in U} \sum_{x^+ \in X} P(x^+ | x, u) \cdot \left( \mathcal{L}(x, u) + \gamma \cdot v(x^+) \right)$

$T_\pi : \mathbb{R}^X \to \mathbb{R}^X$   — affine

$: v \longmapsto (T_\pi v)(x) = \sum_{u \in U} \pi(u | x) \cdot \sum_{x^+ \in X} P(x^+ | x, u) \cdot \left( \mathcal{L}(x, u) + \gamma \cdot v(x^+) \right)$

$$\ddots : v \longmapsto (T_\pi v)(x) = \sum_{u \in \mathcal{U}} \pi(u|x) \cdot \sum_{x^+ \in X} P(x^+|x,u) \cdot \left( \mathcal{L}(x,u) + \gamma \cdot v(x^+) \right)$$

$*$ recall that $T \& T_\pi$ are contractions: $\quad \|Tv - Tw\|_\infty \leq \gamma \|v - w\|_\infty$

$$\|T_\pi v - T_\pi w\|_\infty \leq \gamma \|v - w\|_\infty$$

$\longrightarrow$ propose a <u>value iteration</u> algorithm that approximates $v^* \, (= Tv^*)$
(discuss computational complexity)

— starting from any $v \in \mathbb{R}^X$, iteratively evaluate (nonlinear) operator

$\quad T: \quad v \longmapsto Tv \longmapsto T^2 v \longmapsto \cdots \longmapsto T^k v$

$* \quad \|T^k v - v^*\| \leq \gamma^k \|v - v^*\| \quad$ so $\quad \lim_{k \to \infty} T^k v = v^*$

— each iteration (evaluation of $T$) requires $O(|X| \cdot |\mathcal{U}|)$ operations

$\rightsquigarrow$ given $v^*$, can determine an optimal deterministic policy $\pi^* : X \to \mathcal{U}$
by evaluating $T$

$\longrightarrow$ propose a <u>policy iteration</u> algorithm that computes $v^\pi \, (= T_\pi v^\pi)$
and then approximates $\pi^*$ (discuss computational complexity)

— given any $\pi : X \to \Delta(\mathcal{U})$, can compute $v^\pi$ by solving
affine equation $v^\pi = T_\pi v^\pi$ in $O(|X|^3)$ operations
$\qquad\qquad\qquad\qquad\qquad\quad \underset{\text{actually } |X|^{2 \text{ to } 3}}{\uparrow}$

— can improve policy with greedy update:

$$\forall x \in X : \pi^+(x) = \arg\min_{u \in \mathcal{U}} \sum_{x^+ \in X} P(x^+|x,u) \cdot \left( \mathcal{L}(x,u) + \gamma \cdot v^\pi(x^+) \right)$$

$$\forall x \in X : \pi^+(x) = \arg \min_{u \in \mathcal{U}} \sum_{x^+ \in X} P(x^+ | x, u) \cdot \left( \mathcal{L}(x, u) + \gamma \cdot V^\pi(x^+) \right)$$

$*$ it turns out that $\forall x \in X : V^{\pi^+}(x) \leq V^\pi(x)$ and $\exists x^* : V^{\pi^+}(x^*) < V^\pi(x^*)$

<u>and</u>  $\pi \longmapsto \pi^+ \longmapsto \pi^{++} \longmapsto \dots \longmapsto \pi^* \rightsquigarrow V^{\pi^*} = V^*$

i.e. this iteration converges to an optimal (deterministic) policy!

$\rightsquigarrow$ in a <u>finite</u> number of iterations!

$\quad \hookrightarrow$ follows from the fact that $|\mathcal{U}^X| < \infty$, but $|\mathcal{U}^X| = |\mathcal{U}|^{|X|}$