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

refs:

## *Neuro-Dynamic Programming*

Dimitri P. Bertsekas and John N. Tsitsiklis

chapter 5

---

◦ given finite MDP/SOCP $(X, \mathcal{U}, P, c)$ — infinite-horizon, exponentially-discounted

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

$\quad$ s.t. $x^+ \sim P(x,\mathcal{u}), \; |X|, |\mathcal{U}| < \infty \qquad\qquad \hookrightarrow \gamma \in (0,1)$

◦ recall that the value $v^\pi : X \to \mathbb{R}$ of policy $\pi : X \to \Delta(\mathcal{U})$

is defined $\forall x \in X : v^\pi(x_0) = E[c(x,\mathcal{u}) \mid x(0) = x_0, \; \mathcal{u}(t) \sim \pi(x(t))]$

and satisfies Bellman eqn: $= \sum_{\mathcal{u}_0 \in \mathcal{U}} \pi(\mathcal{u}_0 | x_0) \sum_{x_1 \in X} P(x_1 | x_0, \mathcal{u}_0)\left(\mathcal{L}(x_0, \mathcal{u}_0) + \gamma \cdot v^\pi(x_1)\right)$

$\qquad\qquad\qquad = E\left[\mathcal{L}(x_0, \mathcal{u}_0) + \gamma \cdot v^\pi(x_1) \mid x_1 \sim P(x_0, \mathcal{u}_0), \; \mathcal{u}_0 \sim \pi(x_0)\right]$

$$= E\left[ \mathcal{L}(x_0, u_0) + \gamma \cdot v^{\pi}(x_1) \mid x_1 \sim P(x_0, u_0), u_0 \sim \pi(x_0) \right]$$

- so, removing subscripts: $\forall x \in X : v^{\pi}(x) = E\left[ \mathcal{L}(x, u) + \gamma \cdot v^{\pi}(x^+) \right]$
  (& suppressing distributions)

* this form of Bellman equation suggests we can approximate $v^{\pi}$
  from trajectory data using Monte Carlo estimation

---

aside: Monte Carlo estimation (Robbins-Monroe)

- let $w : \Omega \to \mathbb{R}$ be a random variable whose mean we want to estimate

* given samples $\{w^n\}_{n=1}^N$ of this rv then $E[w] \simeq \overline{w}_N := \frac{1}{N} \sum_{n=1}^N w^n$
  $\underbrace{\qquad\qquad}_{\text{sample mean}}$

→ show that the sample mean can be computed recursively
  (find a way to write $\overline{w}_{N+1}$ i.t.o. $\overline{w}_N$ and $w^{N+1}$)

— $\overline{w}_{N+1} = \frac{1}{N+1} \sum_{n=1}^{N+1} w^n = \overline{w}_N + \frac{1}{N+1}\left( w^{N+1} - \overline{w}_N \right)$

- note that, by rearranging, we get $\overline{w}_{N+1} = \left(1 - \frac{1}{N+1}\right)\overline{w}_N + \frac{1}{N+1} w^{N+1}$
  which has the form $x^+ = (1-\alpha)x + \alpha u$

* for any $\alpha \in (0,1)$, if $u$ is rv with mean $E[u]$
  then $x$ is rv with mean $E[x] = E[u]$

○ recalling that $\forall x \in X : v^{\pi}(x_0) = E[c(x,u) \mid x(0) = x_0]$
we can apply Monte Carlo estimation given trajectories $\{(x^n, u^n)\}_{n=1}^{N+1}$
generated by applying policy $\pi$ from initial condition $x^n(0) = x_0$:

$$v_N^{\pi}(x_0) = \frac{1}{N} \sum_{n=1}^{N} c(x^n, u^n), \quad v_{N+1}^{\pi}(x_0) = v_N^{\pi}(x_0) + \alpha \cdot \left( c(x^{N+1}, u^{N+1}) - v_N^{\pi}(x_0) \right)$$

stylize

$\alpha \in (0,1), \text{ eg } \alpha = \frac{1}{N+1}$

$$v^+(x_0) = v(x_0) + \alpha \left( c(x,u) - v(x_0) \right), \quad c(x,u) = \sum_{t=0}^{\infty} \gamma^t \cdot \mathcal{L}(x_t, u_t)$$

$$= v(x_0) + \alpha \cdot \Big[ \gamma^0 \cdot \big( \mathcal{L}(x_0, u_0) + \gamma v(x_1) - v(x_0) \big)$$

$$+ \gamma^1 \cdot \big( \mathcal{L}(x_1, u_1) + \gamma \cdot v(x_2) - v(x_1) \big)$$

$$+ \cdots$$

$$+ \gamma^t \cdot \big( \mathcal{L}(x_t, u_t) + \gamma \cdot v(x_{t+1}) - v(x_t) \big)$$

$$+ \cdots \Big]$$

$$* \quad v^+(x_0) = v(x_0) + \alpha \cdot \sum_{t=0}^{\infty} \gamma^t \cdot \underbrace{\big( \mathcal{L}(x_t, u_t) + \gamma \cdot v(x_{t+1}) - v(x_t) \big)}$$

$=: d_t$ — temporal differences (TD)

$\hookrightarrow$ this is TD(0) — generalizes to TD($\lambda$)

Q-learning

○ given policy $\pi : X \to \Delta(u)$ with value $v^{\pi} : X \to \mathbb{R}$,
consider state-action policy quality

$\begin{cases} 1^\circ: \text{choosing } u \text{ in } x \\ 2^\circ: \text{using } \pi \text{ from then on} \end{cases}$

$$q^{\pi} : X \times U \to \mathbb{R} \quad - \text{ expected cost of :}$$

$$: (x, u) \mapsto \sum_{x^+ \in X} P(x^+ \mid x, u) \cdot \left( \mathcal{L}(x,u) + \gamma \cdot v^{\pi}(x^+) \right)$$

○ note that optimal quality function satisfies the Bellman equation

$x \in X$

• note that optimal quality function satisfies the Bellman equation

$$Q^*(x,u) = E\left[ \mathcal{L}(x,u) + \gamma \cdot \min_{w \in \mathcal{U}} Q^*(x,w) \right] \quad b/c \quad v^k(x) = \min_{u \in \mathcal{U}} Q^*(x,u)$$

* applying Monte-Carlo estimation yields $\underline{Q\text{-learning}}$ algorithm

$$Q^+(x,u) = (1-\alpha) Q(x,u) + \alpha \cdot \left( \mathcal{L}(x,u) + \gamma \cdot \min_{w \in \mathcal{U}} Q(x,w) \right)$$

$$\uparrow \; \alpha \in (0,1) - \text{learning rate}$$