# Nonlinear function approx$^n$ in RL

Structural conditions:

$Q^*$ realizable $\qquad Q^* \in \mathcal{F} / Q$

local optimism
$$Q_h^*(s,a) \leq Q_h(s,a) \quad \forall \, h,s,a$$

linear MDP - low rank enough & LVI-UCB works

$\underline{P/R}$ are linear

$\Downarrow$

$Q \leftarrow$ linear + bonus

$\underline{\qquad\qquad}$
linear

(linear) Bellman completeness.

$$\underline{Q \in \mathcal{F}} \qquad \Rightarrow \qquad TQ \in \mathcal{F} \qquad \leftarrow$$

How to measure complexity of $\mathcal{F}$ for RL?

Finite-horizon MDP $\quad M = (S, A, P, R, \mu, H)$

Realizability $\quad Q_h^*(s,a) = \langle \theta_h^*, \phi(s,a) \rangle \quad$ for linear $\quad \in \mathcal{F}$ for nonlinear

Bellman completeness (linear) For any $\theta$, $\exists\, \bar{\theta}$ s.t. $\qquad \theta \equiv Q$
$$(T\theta)(s,a) = \langle \bar{\theta}, \phi(s,a) \rangle$$

Global optimistic regret decomposition $\quad$ If $Q_1 \ldots Q_H$ s.t.

$\Rrightarrow E_{s_1} \max\limits_a Q_1(s_1,a) \geqslant E_{s_1} \max\limits_a Q_1^*(s_1,a) \qquad$ (earlier $Q \geqslant Q^*$ optimistic

$\forall s,a,h$

and $\pi$ greedy w.r.t. $Q$, then $\qquad Q \leq TQ + \text{conf}$

$$J(\pi^*) - J(\pi) \leq \sum_{h=1}^{H} E_{(s,a) \sim d_h^\pi}[Q_h(s,a) - (TQ_{h+1})(s,a)]$$
$$E[V^*(s) - V_\pi(s_1)] \qquad\qquad \underbrace{\qquad\qquad}_{\text{conf}_h(s,a)}$$

Proof: $\quad J(\pi^*) - J(\pi) = E_s[\underbrace{Q^*(s_1, \pi^*(s_1))}_{\max\limits_a Q^*(s_1,a)} - Q_1^\pi(s_1, \pi(s_1))]$

by
global $\quad \leq E_{s_1}[\max\limits_a Q_1(s_1,a) - Q_1^\pi(s_1, \pi(s_1))] \quad ①$
opt

$\pi$ is greedy $\leq E_{(s_1,a_1) \sim d_0^\pi}[Q_1(s_1, a_1) - Q_1^\pi(s_1, \pi(s_1))]$
w.r.t. $Q_1 \qquad\qquad\qquad\qquad \underset{\pi(s_1)}{\,}$

$= E_{(s,a) \sim d_0^\pi}[Q_1(s_1,a_1) - \underline{(TQ_2)(s_1,a_1)}$
$\qquad\qquad + \underline{(TQ_2)(s_1,a_1)} - Q_1^\pi(s_1, \pi(s_1))]$

$$= \text{first term} + R(s_1, a_1) + E_{s_2}[\max_a Q_2(s_2, a_2)]$$

$$\underbrace{Q_1 - TQ_2}_{} \qquad - R(s_1, a_1) - E_{s_2}[Q_2^\pi(s_2, \pi(s_2))] \quad \text{②}$$

$$= \sum_{h=1}^{H} E[Q_h - TQ_{h+1}] \qquad \underbrace{\qquad\qquad}_{\text{B}}$$

**Main challenge:** Show global optimism only assuming Bellman completeness.

[Linear] Bellman completeness : For any $\underline{\theta}$, $\exists \bar{\theta}$ st.

$$\forall (s,a) \qquad (T\theta)(s,a) = \langle \bar{\theta}, \phi(s,a) \rangle$$

$$\| $$

$$E[R(s,a) + \max_{a'} Q(s', a')]$$

$$\| $$

$$\langle \theta, \phi(s', a') \rangle$$

$\Rightarrow$ no misspecification in regression. if take $Q$ to be linear.

① Regression error $R_h^{t-1}(\theta, \tilde{\theta}_{h+1}) = \sum_{i=1}^{t-1} \left( \langle \phi(s_h^i, a_h^i), \theta \rangle - r_h^i - \underbrace{\max_a \langle \phi(s_{h+1}^i, a), \tilde{\theta}_{h+1} \rangle}_{V_{h+1}^i(s)} \right)^2 \leftarrow$

$$\underline{\underline{\qquad}} \qquad \equiv (Q - T\theta)^2 \qquad \underbrace{\qquad\qquad\qquad}_{Y_i}$$

earlier $Q$ = linear + bonus

how $Q$ = linear

no optimism!

earlier $EY_i \neq$ linear.

now $EY_i$ = linear acc to Bellman completeness.

Build in optimism using nested confidence balls

② $BALL^t := \{(\theta_1, ..., \theta_H): \theta_H = 0, \forall h \; R_h^{t-1}(\theta_h, \theta_{h+1}) \leq \min_\theta R_h^{t-1}(\theta, \theta_{h+1}) + \beta^2\}$

$\qquad \underset{=}{} \qquad\qquad \underset{\neq}{} $ optimism

$\max_{\theta \in C_x} \chi\theta$
$\theta^* \in C_x$

Note that $\theta^*$ : $R_h^{t-1}(\theta_h^*, \theta_{h+1}^*) = 0$

$$Q^* = TQ^*$$

$$\theta_h^* = R(s_{h,L}) + E[V_{h+1}^*(s')]$$

$SQ_h^*\}_{h \in H} \equiv \theta_1^* - \theta_H^*$

③ $\Rightarrow$ $(\theta_1^t, ..., \theta_H^t) \leftarrow$ argmax $E_{s_1} \max_a \langle \phi(s, a_1), \theta_1 \rangle$
$(\theta_1, ... \theta_H) \in BALL^t$

$(\theta_1^*, ... \theta_H^*) \in BALL^t$

$\qquad\qquad\qquad E_{s_1} \max \langle \phi(s_1, a_1), \theta_1^* \rangle$

$$\Rightarrow \text{Global optimism holds} \quad \text{...} \quad a$$
$$\leq E_{s_1} \max_a \langle \phi(s_1,a_1), \hat{\theta}_1 \rangle$$

⑥ greedy policy wrt $\theta^t$ to collect another episode.

Regret analysis:

$$\text{Regret} \leq E\left[\sum_t \sum_h \theta_h^t(s_h^t, a_h^t) - (T\theta_{h+1}^t)(s_h^t, a_h^t)\right] \equiv E\left[\sum_t \sum_h \text{conf}_h^t\right]$$

$$\leq \sum_t \sum_h \theta_h^t(s_h^t, a_h^t) - (T\theta_{h+1}^t)(s_h^t, a_h^t) + \tilde{O}(H\sqrt{T})$$

depend on
linear
assumption

$$\leq \sum_t \sum_h \langle \phi(s_h^t, a_h^t), \theta_h^t - \bar{\theta}_h^t \rangle + \tilde{O}(H\sqrt{T})$$

$$\leq \sum_t \sum_h \|\phi(s_h^t, a_h^t)\|_{\Lambda_{B,t-1}^{-1}} \cdot \|\theta_h^t - \bar{\theta}_h^t\|_{\Lambda_{B,t-1}} + \tilde{O}(H\sqrt{T})$$

$\underbrace{\qquad\qquad\qquad}_{\text{elliptic potential lemma}}$ $\underbrace{\qquad\qquad}_{\beta = O(H\sqrt{d})}$ saving $\sqrt{d}$ factor using global opt.

$$= \tilde{O}(H\beta\sqrt{dT}) + \tilde{O}(H\sqrt{T})$$

$$= \tilde{O}(H^2 d\sqrt{T})$$ ▨

Generalization to nonlinear functions.
 - only place linear assumption needed is to bound $\|\theta_h^t - \bar{\theta}_h^t\|$
   $\underbrace{\qquad}_{d}$

Notion for nonlinear complexity

Bellman rank : Given $F$, let $\Pi$ be induced policy class

$\Pi = \{\pi_f : f \in F\}$. For each $h$, $\exists$ embedding function
  $w_h : \Pi \to \mathbb{R}^d$ and $v_h : F \to \mathbb{R}^d$ s.t.

Bellman error $\mathcal{E}_h(\Pi, F) = \langle w_h(\pi), v_h(f) \rangle$    Q~TQ

where $\mathcal{E}_h(\Pi, F) = E[\theta_h(s_h, a_h) - r_h - \max_a \theta_{h+1}(s_{h+1}, a')|$

$$s_h \sim d_h^\pi, a_h = \pi_\theta(s_h)]$$
not $\pi$ as for linear

d- Bellman rank

⌈Note: $(s \neq \text{linear})$

Regret:

$$J(\pi^*) - J(\pi) = \sum_t \sum_h \langle \omega_h(\pi^t), v_h(f^t) \rangle \qquad \left( \substack{\text{earlier} \\ \sum_t (Q - T\theta)} \right)$$

$$\leq \sum_t \sum_h \| \omega_h(\pi^t) \|_{\Sigma_{t-1,h}^{-1}} \| v_h(f^t) \|_{\Sigma_{t+1,h}}$$

$$\subset \lambda I + \sum v v^\top$$

assuming $\| \omega_h(\pi) \|_2 \leq W$, $\| v_h(f) \|_2 \leq V$

can show $\| v_h(f^t) \|_{\Sigma_{t+1,h}} \leq \sqrt{\lambda V^2 + 4\beta^2} = \| \theta - \hat{\theta} \|$

$$\| \omega_h(f^t) \|_{\Sigma_{t+1,h}^{-1}} \leq \sqrt{\frac{2Hd}{T} \log |\cdot|} \quad ) \quad \substack{\text{elliptic} \\ \text{potential} \\ \text{lemma}}$$

$\Rightarrow$ Regret bound that depends on Bellman rank $d$.