

Linear MDP A linear MDP $M = (S, A, P, R, H, d_0)$
 satisfies $P(s'|s, a) = \phi(s, a)^T \psi(s')$ ($\phi, \psi \in \mathbb{R}^d$, $d \ll |S \times A|$)
 $R(s, a) = \phi(s, a)^T \theta_R$. "low-rank MDP"
 learner knows $\phi: S \times A \rightarrow \mathbb{R}^d$. $R_{\max} = 1$.

Key property: $\forall f: S \times A \rightarrow \mathbb{R}, T f, T^\pi f \in \mathcal{F}$.
 $\mathcal{F} := \{ (s, a) \mapsto \phi(s, a)^T \theta : \theta \in \mathbb{R}^d \}$.
 $(Tf)(s, a) = R(s, a) + \langle P(s, a), V_f \rangle$
 $\Delta = \phi^T(s, a) \theta_R + \phi(s, a)^T (\Psi V_f)$
 $= \phi^T(s, a) (\theta_R + \Psi V_f)$
 $\Rightarrow Q^*, Q^\pi \in \mathcal{F}$.

Known vs. unknown. $\mathbb{R}^{d \times d}$. \mathbb{R}^d
 $\forall h \in \{1, 2, \dots, H\}, \Lambda_h = \mathbb{E}_{D_h} [\phi(s_h, a_h) \phi(s_h, a_h)^T]$
 $\{ s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_H, a_H, r_H \}$

"Known" (s_n, a_n) : when $\frac{\phi(s_n, a_n)^T \Lambda_n^{-1} \phi(s_n, a_n)}{\Delta}$ is small
 "unknown" (s_n, a_n) : if Δ is big.

Special case: tabular: $\phi(s, a)$ is 1-hot.

Λ_n is diagonal with $\Lambda_n(s, a), (s, a) = \frac{n(s, a)}{n}$.

Alg: at each episode,

$$f_{H+1} \equiv 0.$$

for $h = H, H-1, H-2, \dots, 1$.

"Optimistic FQI"

$$\hat{f}_h = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \mathbb{E}_{\mathcal{D}_h} [(f(s, a) - r - \underbrace{V_{f_{h+1}}(s')})^2]$$

$$f_h(s, a) = \begin{cases} \hat{f}_h(s, a), & \text{if } (s, a) \in K. \\ V_{\max}, & \text{if } (s, a) \notin K. \end{cases}$$

\hookrightarrow non-smooth.

$$f := f_1 \circ f_2 \circ f_3 \circ \dots \circ f_H.$$

explore w/ $\pi_f := \pi_{f_1} \circ \pi_{f_2} \circ \dots \circ \pi_{f_H}$.
 (collect a batch sample).

Simplification: batch size $\rightarrow \infty$. (only polynomial)

want to bound iteration complexity.

$$\Rightarrow \begin{cases} (s,a) \in K: \phi(s,a)^T \Lambda_h^{-1}(s,a) \phi(s,a) < +\infty \\ (s,a) \notin K: \phi(s,a)^T \Lambda_h^{-1}(s,a) \phi(s,a) = +\infty \end{cases} \checkmark$$

will show: in each round, $\exists h, s.t.$
 $\text{rank}(\Lambda_h)$ inc by at least 1.

$$V_{\max} = (H-h+1) \cdot R_{\max}$$

Lemma. (Optimism): $f \geq Q^*$

$$\text{At } H: f_H(s,a) = \begin{cases} R(s,a) = Q_H^*(s,a) \cdot \forall (s,a) \in K \\ V_{\max} \geq Q_H^*(s,a) \end{cases}$$

For $h < H$, by induction $f_{h+1} \geq Q_{h+1}^*$

$$\Rightarrow V_{f_{h+1}} \geq V_{h+1}^*$$

$$f_h(s,a) = \begin{cases} \mathbb{E}[r + V_{f_{h+1}}(s') | s,a] \geq \mathbb{E}[r + V_{h+1}^*(s') | s,a] \\ \phantom{\mathbb{E}[r + V_{f_{h+1}}(s') | s,a]} = Q_h^*(s,a) \end{cases}$$

\downarrow rely on
 $* \mathbb{E}[\dots | s,a] = (T f_{h+1})(s,a)$
 $* \forall f, T f \in \mathcal{F}$

"Optimal - or - explore"

$$\begin{aligned}\epsilon &< J(\pi^*) - J(\pi_f) = \mathbb{E}_{s \sim d_0} \left[\max_a Q^*(s, a) \right] \\ &\leq \mathbb{E}_{s \sim d_0} \left[\max_a f_1(s, a) \right] - J(\pi_f) \\ &= \mathbb{E}_{s \sim d_0} \left[\underbrace{f(s, \pi_f)}_{\text{viewed as prod of } J(\pi_f)} \right] - J(\pi_f).\end{aligned}$$

$$= \sum_{h=1}^H \mathbb{E}_{\pi_f} \left[f(s_h, a_h) - r_h - V_f(s_{h+1}) \right].$$

$$\Rightarrow \exists h, \frac{\epsilon}{H} \leq \mathbb{E}_{\pi_f} \left[f(s_h, a_h) - r_h - V_f(s_{h+1}) \right]$$

✓ $f \geq \tilde{f}$

$$\leq \mathbb{E}_{\pi_f} \left[f(s_h, a_h) - r_h - V_{\tilde{f}}(s_{h+1}) \right].$$

$$\leq \left| \mathbb{E}_{\pi_f} \left[f(s_h, a_h) - \mathbb{E} \left[r_h + V_f(s_{h+1}) \mid s_h, a_h \right] \right] \right|$$

$$\leq \sqrt{\mathbb{E}_{\pi_f} \left[\left(f(s_h, a_h) - \mathbb{E} \left[r_h + V_f(s_{h+1}) \mid s_h, a_h \right] \right)^2 \right]}$$

(*)

$$\mathbb{E}_{\pi_f}[\dots] = \sum_{(s,a)} d_n^{\pi_f}(s,a) \cdot (\dots)$$

$$= \sum_{\substack{(s,a) \in K \\ \text{w.t.s.} \\ \text{zero.}}} d_n^{\pi_f}(s,a) (\dots)$$

$$+ \sum_{(s,a) \notin K} d_n^{\pi_f}(s,a) (\dots)$$

$$\sum_{(s,a) \in K} d_n^{\pi_f}(s,a) \left(f(s,a) - \mathbb{E}[r + V_f(s') | s,a] \right)^2$$

$$\forall (s,a) \in K$$

$$f_n(s,a) = \tilde{f}_n(s,a) = \mathbb{E}[r + V_f(s') | s,a]$$

needs to be shown.

Linear regression. design $X \in \mathbb{R}^{n \times d}$ output $Y \in \mathbb{R}^{n \times 1}$.

$$\underline{x^T \left((X^T X)^{-1} X^T Y \right)}$$

with $y = x^T \theta + \varepsilon$
 $\varepsilon = \text{zero mean.}$

$$\underline{x^T \left((X^T X)^{-1} X^T Y \right)}$$

$$x_0 \in \mathbb{R}^{d \times 1}$$
$$y_0 = x_0^T \theta$$

$$a x_0^T \underbrace{(x_0 x_0^T)^{-1}} x_0^T y_0$$
$$= a y_0.$$

$$\sum_{(s,a) \in K} d^{af}(s,a) \underbrace{(\dots)}_{\leq O(V_{\max})} \geq \frac{\varepsilon}{H}.$$

$$\Rightarrow \sum_{(s,a) \in K} d^{af}(s,a) \geq \frac{\varepsilon}{H V_{\max}} > 0.$$

If we explore w/ π_f .

will collect data from $(s, a) \notin K$.

$$\Lambda'_n = \alpha \Lambda_n + (1 - \alpha) \mathbb{E}_{\pi_f} [\underbrace{\phi(s, a) \phi(s, a)^T}_{\text{rank } 1}]$$

Suppose $x^T \Lambda^{-1} x = +\infty$
 w.t.s, that $\Rightarrow d \times d$.

$$\text{rank}(\Lambda + \alpha \alpha^T) > \text{rank}(\Lambda)$$

$\phi(s, a) \notin K$.

\downarrow

$$\phi(s, a)^T \Lambda_n^{-1} \phi(s, a) = +\infty$$

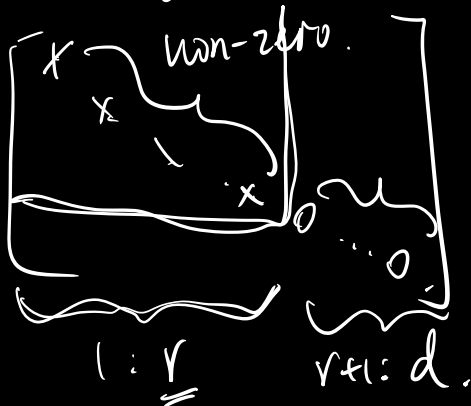
Proof Sketch: eigen-decompose

$$\Lambda = U^T \Sigma U \in \mathbb{R}^{d \times d}$$

$$\tilde{x} = U x$$

$$x^T \Lambda^{-1} x = x^T U^T \Sigma^{-1} U x$$

$$= \tilde{x}^T \Sigma^{-1} \tilde{x} = +\infty$$

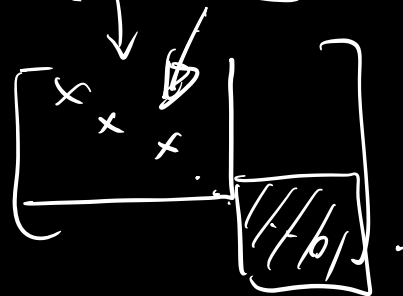


$$\Rightarrow [\tilde{x}]_{r+1:d} \neq \vec{0}$$

$$\text{rank}(\Lambda + \alpha \alpha^T)$$

$$= \text{rank}(U^T \Sigma U + U^T \tilde{x} \tilde{x}^T U)$$

$$= \text{rank}(\Sigma + \tilde{x} \tilde{x}^T)$$



Remarks:

1. If we have explored everywhere,
i.e. $\mathbb{E}[\phi\phi^\top]$ is full-rank,

(Optimistic) FQI must output a near-optimal.

but in FQI lecture. we define exploratory data μ to be such that $\max_{\pi} \frac{d_{\pi}^{\pi}(s,a)}{\mu(s,a)} \leq C$.

Can unify: in FQI, we need C thru:

$$\sup_{f, f' \in \mathcal{F}} \frac{\|f - f'\|_{2, d_{\pi}^{\pi}}}{\|f - f'\|_{2, \mu}} \leq \frac{d_{\pi}^{\pi}(s,a)}{\mu(s,a)} \leq C.$$

$$f - Tf, \quad Tf \in \mathcal{F}.$$

✓

can show: when \mathcal{F} is linear,

⊗ $\mathbb{E}_{\mu}[\phi\phi^\top]$ has full rank

this quantity $< +\infty$.

2. UCB-LSVI [Jin et al '20].

→ $Tf \in \mathcal{F} \quad \forall f$

→ use non-smooth.

Question: can we do exploration with (say).

$Q^* \in \mathcal{F}$?

In low-rank MDP. (ie. linear MDP w/o Φ given to learner)

$Q^* \in \mathcal{F} \Rightarrow ?$

Bellman rank [JKALS'17].

$$1. \quad X \in \mathbb{R}^{n \times d}, \quad Y \in \mathbb{R}^{n \times 1}, \quad Y = X\theta, \quad \Lambda = X^T X.$$

Give $x_0 \in \mathbb{R}^d$ s.t. $x_0^T \Lambda^{-1} x_0 < +\infty$. $\lim_{\lambda \rightarrow 0} (\Lambda + \lambda I)^{-1}$

show that we can recover $x_0^T \theta$ even when Λ is rank deficient

Proof: eigen: $\Lambda = U^T \Sigma U$ $(Ux_0)^T \Sigma^{-1} (Ux_0) < +\infty$.
 w.l.o.g. assume $\Sigma = \begin{bmatrix} x_1 & & & \\ & \ddots & & \\ & & x_r & \\ & & \underbrace{0 \dots 0}_{\text{non-zero}} & \\ & & & \ddots & \\ & & & & 0 \end{bmatrix}$ ($r = \text{rank}(\Lambda)$)
 $1: r \quad r+1: d$

$$\Rightarrow [Ux_0]_{r+1:d} = \vec{0}$$

$$\text{v.t.s.} \quad x_0^T \theta = x_0^T (\Lambda^{-1} X^T Y)$$

$$\text{RHS} = x_0^T \Lambda^{-1} X^T Y \quad x_0^T \Lambda^{-1} X^T$$

$$= x_0^T \Lambda^{-1} X^T X \theta$$

$$= x_0^T \Lambda^{-1} \Lambda \theta$$

$$= x_0^T U^T \Sigma^{-1} U U^T \Sigma U \theta$$

$$= x_0^T U^T \Sigma^{-1} \Sigma U \theta.$$

$$= (Ux_0)^T \Sigma^{-1} \Sigma U \theta = (Ux_0)^T I_d U \theta$$

$$\begin{bmatrix} x_1 \\ \vdots \\ x_r \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & 0 \dots 0 \end{bmatrix}$$

$1: r$

$$= x_0^T U^T U \theta$$

$$= x_0^T \theta \quad \checkmark$$

2. Given $x_1^T \Omega^{-1} x_1 = +\infty$.

show $\text{rank}(\Omega + x_1 x_1^T) > \text{rank}(\Omega)$.

(and similarly, if $x_1^T \Omega^{-1} x_1 < +\infty$, rank doesn't change)

Proof: similar to above, we do everything in the eigen-decomposition of Ω , so.

w.l.o.g., we can assume Ω is diagonal.

$$\Omega = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & \ddots \\ & & & & 0 & \dots & 0 \\ & & & & \vdots & \dots & \vdots \\ & & & & & & 0 \end{bmatrix} \quad (\text{also normalize the diagonals to } 1)$$

$1:r \quad r+1:d$

$$x_1^T \Omega^{-1} x_1 = \infty \iff [x_1]_{r+1:d} \neq \vec{0}$$

$$\text{Let } x_1 = \begin{bmatrix} \bar{x}_1 & \tilde{x}_1 \\ 1:r & r+1:d \end{bmatrix}$$

$$\Rightarrow \Omega + x_1 x_1^T = \begin{bmatrix} I_r + \bar{x}_1 \bar{x}_1^T & \bar{x}_1 \tilde{x}_1^T \\ \tilde{x}_1 \bar{x}_1^T & \tilde{x}_1 \tilde{x}_1^T \end{bmatrix}$$

now try to block-diagonalize this matrix.

$$\begin{bmatrix} I_r + \bar{x}_1 \bar{x}_1^T & \bar{x}_1 \tilde{x}_1^T \\ \tilde{x}_1 \bar{x}_1^T & \tilde{x}_1 \tilde{x}_1^T \end{bmatrix} = \begin{bmatrix} I_r & 0 \\ v & I_{d-r} \end{bmatrix} \begin{bmatrix} I_r + \bar{x}_1 \bar{x}_1^T & 0 \\ 0 & z \end{bmatrix} \begin{bmatrix} \Sigma & v^T \\ 0 & I \end{bmatrix}$$

w.t.s. $z \neq 0$

$$\begin{bmatrix} (I_r + \bar{x}_1 \bar{x}_1^T) & 0 \\ V(I_r + \bar{x}_1 \bar{x}_1^T) & Z \end{bmatrix} \begin{bmatrix} I & V^T \\ 0 & I \end{bmatrix}$$

$$= \begin{bmatrix} I_r + \bar{x}_1 \bar{x}_1^T & (I_r + \bar{x}_1 \bar{x}_1^T) V^T \\ V(I_r + \bar{x}_1 \bar{x}_1^T) & V(I_r + \bar{x}_1 \bar{x}_1^T) V^T + Z \end{bmatrix} \stackrel{\text{rank}}{=} \begin{bmatrix} I_r + \bar{x}_1 \bar{x}_1^T & \tilde{x}_1 \tilde{x}_1^T \\ \tilde{x}_1 \bar{x}_1^T & \tilde{x}_1 \tilde{x}_1^T \end{bmatrix}$$

$$\Rightarrow (I_r + \bar{x}_1 \bar{x}_1^T) V^T = \tilde{x}_1 \tilde{x}_1^T$$

$$V^T = (I_r + \bar{x}_1 \bar{x}_1^T)^{-1} (\tilde{x}_1 \tilde{x}_1^T)$$

$$\Rightarrow Z = \tilde{x}_1 \tilde{x}_1^T - \underbrace{(\tilde{x}_1 \bar{x}_1^T)}_{(d-r) \times r} \underbrace{(I_r + \bar{x}_1 \bar{x}_1^T)^{-1}}_{r \times r} \underbrace{(\tilde{x}_1 \tilde{x}_1^T)}_{r \times d}$$

to show $Z \neq 0$,

suffices to show $\exists y, y^T Z y > 0$.

$$\text{let } y = \begin{pmatrix} \hat{\tilde{x}}_1 \\ \tilde{x}_1 \end{pmatrix} \Rightarrow y^T \tilde{x}_1 \tilde{x}_1^T y = \|\tilde{x}_1\|^2$$

$$y^T \tilde{x}_1 \underbrace{(\bar{x}_1^T (I_r + \bar{x}_1 \bar{x}_1^T)^{-1} \bar{x}_1)}_{\text{w.t.s.} < 1} \tilde{x}_1^T y$$

$$\underbrace{y^T \tilde{x}_1}_{\|\tilde{x}_1\|}$$

w.t.s. < 1 .

$$\|\tilde{x}_1\|$$

$$\bar{x}_1^T (I_r + \bar{x}_1 \bar{x}_1^T)^{-1} \bar{x}_1$$

rotate into basis
w/ (\bar{x}_1) as a basis vector

let Q be orthogonal s.t.

$$\bar{x}_1 = Q \begin{bmatrix} \|\bar{x}_1\| \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\begin{aligned} & [\|\bar{x}_1\|, 0, \dots, 0] Q^T (I_r + Q \begin{bmatrix} \|\bar{x}_1\|^2 \\ 0 \\ \vdots \\ 0 \end{bmatrix} Q^T)^{-1} Q \begin{bmatrix} \|\bar{x}_1\| \\ 0 \\ \vdots \\ 0 \end{bmatrix} \\ &= \frac{\|\bar{x}_1\|^2}{\|\bar{x}_1\|^2 + 1} < 1 \end{aligned}$$