

Offline learning

Infinite-horizon

$s_0 \rightarrow$  lower bound only.

$M = (S, A, P, R, \gamma, d)$   $d \in M_1(S)$  initial state dist.

$$V_M^\pi(s) = \mathbb{E}^{M, \pi} \left[ \sum_{h=0}^{\infty} \gamma^h r_h \mid S_0 = s \right] ; Q_M^\pi(s, a) = \mathbb{E}^{M, \pi} \left[ \sum_{h=0}^{\infty} \gamma^h r_h \mid S_0 = s, A_0 = a \right]$$

$\exists$  determ. & ML policy  $\pi_M^*$  s.t.  $V_M^* = V_M^{\pi_M^*} = \max_{\pi} V_M^\pi$ ,  $Q_M^* = Q_M^{\pi_M^*}$ ,  $\pi_M^* = \operatorname{argmax}_{\pi} Q_M^\pi$

Dataset

oracle-generated  $D = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^n$   $(s_i, a_i) \sim M \in M_1(S \times A)$

Policy-induced  $(T_{\log}, \vec{h})$

$D = \{t_j\}_{j=1}^m$   $t_j = (s_0^{(j)}, a_0^{(j)}, r_0^{(j)}, \dots, s_{n_j-1}^{(j)}, a_{n_j-1}^{(j)}, r_{n_j}^{(j)}, s_{n_j}^{(j)})$

Goal [Data is collected before any interactions with the env.]

Batch policy opt.  $(s_0, \{M \in M\})$  share the same state, action &  $\gamma$

w.p.  $> 1 - \delta$ , for every instance  $(s_0, M)$ , when deciding the collection dist., we cannot dist.

the learnt policy  $\hat{\pi}_D$  is near opt.  $\text{IP}(|V^*(s_0) - V^{\hat{\pi}_D}(s_0)| \leq \varepsilon) \geq 1 - \delta$

$$\forall w, \min_w \mu(C_\delta(w) \cup C_{\delta+t}(w)) \leq \frac{2\text{Vol}(C_\delta(w))}{\text{Vol}(B)} = \frac{2}{N(\delta)}$$

$\oplus P_w(\{a_1, \dots, a_N\} \notin C_\delta(w) \cup C_{\delta+t}(w))$

$$\geq (1 - \frac{2}{N(\delta, d)})^n \geq (1 - \frac{2}{N(\delta, d)})^{\frac{N}{d}-1} \geq \frac{1}{e}$$

$$\text{if } n < \frac{N(\delta, d)}{2}/2$$

## Assumptions:

1. realizability (remove dependence from  $|S|$  &  $|A|$ )

Given an instance  $(s_0, M)$ ,  $\exists$  feat. map  $\phi(s, a) \in \mathbb{R}^d$  with  $\|\phi(s, a)\|_2 \leq 1$

such that  $\forall M \in \mathcal{M}$  and  $\forall \pi$ ,  $\Theta(s, a) = \phi(s, a)^T \theta_M^\top$  with  $\|\theta_M^\top\|_2 \leq 1$

Thm 3:  $\exists$  a BPO  $(s_0, M)$  with all-policy realizability s.t.

any algo. using less than  $\frac{N(\epsilon, d)}{2}$  samples has

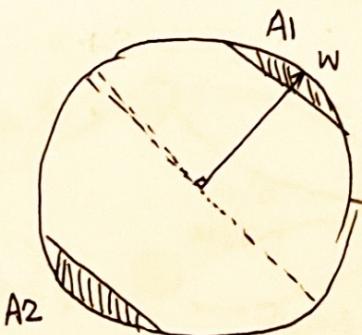
$$V^*(s_0) - V^{\pi_\theta}(s_0) \geq 1 \text{ w.p. at least } 1 - \frac{1}{2e} \quad (\text{Zanette 2020})$$

B

① Rewards are not seen.

② Bootstrapping issue, or the ~~rank~~ deficient features cause the hardness.

A Proof steps: 1. realizability holds 2. the err. of not choosing the opt. action is large 3.  $\exists$  two MDPs which cannot be distinguished under  $n$  samples but has diff. opt. actions



D

- ① rewards only exist in small areas  $\{x: x^T w \geq \sigma\}$
- ② unshadowed area:  $(\phi(s) - \sigma \phi(s'))^T w = 0 \leq -\sigma$   
 $\Rightarrow$  there are multiple possible soln. without exploring the small caps

C  $B = \{x: \|x\|_2 \leq 1\}$

$$\leq \frac{2}{(\frac{1-\gamma}{1-\gamma})^d}$$

$$S = \{s_0\} \cup B$$

$$A(s) = \begin{cases} B, & \text{if } s = s_0 \\ S, & \text{o.w.} \end{cases}$$

$$\phi(s, a) = a$$

$$s'(a) = \begin{cases} \frac{1}{\sigma} (a^T w) \cdot w, & \text{if } a \notin G(w) \cup G(-w) \\ a, & \text{if } a \in G(w) \cup G(-w) \end{cases}$$

③  $r_w = 0$  in the unshadowed area.

$+w, -w$  can be soln.

$M_+$   $M_-$

{ positive rewards if  $a^T w > 0$

neg. o.w.

For  $M_+$ , opt. arm  $+w$  with  $q_+(s_0, w) = 1$

$$q_-(s_0, w) = -1$$

## 2. Concentrability (a) feature coverage

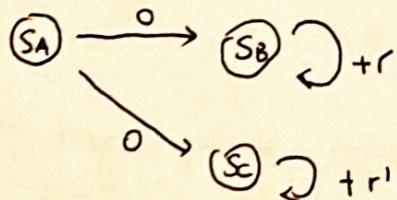
$\theta(s, a)$ ,  $\|\psi(s, a)\|_2 \leq 1$  and  $E_{(s, a) \sim \mu} [\psi(s, a) \psi(s, a)^T]$  has its smallest eigenvalue equaling to  $\frac{1}{d}$ .

how feat. cover all the dimensions

supervised learning

largest var. dir.  $\Rightarrow$  eig. vec.  
its magnitude  $\Rightarrow$  eig. val.

(Amortila 2020)



$$\begin{aligned}\psi(S_A, a) &= [0, \gamma]^T, \quad \psi(S_C, a) = [0, 1]^T \\ \psi(S_B) &= [\gamma, 0]^T \\ D &= \{(S_A, 0, S_B)\}\end{aligned}$$

$$\theta^* = \frac{\gamma}{1-\gamma}$$

$$E_{S \sim M} [\phi(s) \phi(s)^T] = \gamma^2 I$$

$q^\pi$  realizability but  $q^\pi$  is not observed ; only  $r + \gamma \hat{v}(s')$  is observed

[linear Bellman realizability]  
closedness

(b) concentrability :  $v \in \mathcal{M}(S \times A)$  admissible, if  $\exists \pi$  &  $n$  s.t.  $\mathbb{P}^\pi(S_n=s, A_n=a) = v(s, a)$

$\exists$  a const.  $C_{\text{conc}} < \infty$  s.t.  $\theta$  admissible  $v$ ,  $\sup_{(S, a)} \frac{|v(S, a)|}{w(S, a)} \leq C$

Thm :  $\forall S \geq q$  and  $\gamma \in (\frac{1}{2}, 1)$ ,  $\exists M$  with  $|S| = S$ ,  $|A|=2$ ,  $|f|=2$  and a sampling dist.  $M$  which is admissible

1.  $q^\pi$  realizable &  $C_{\text{conc}} \leq C$

2.  $A \leftarrow e - S^{\frac{1}{2}}$  gives  $v^\pi(S_0) =$

2.  $n \leq c \cdot \min \{S^{\frac{1}{3}}(\log S)^2, 2^{C/32}, 2^{\frac{1}{1-\gamma}}\}$  must have  $D(\pi^*) - E_{D, M}[\hat{D}(\hat{\pi})] \geq C'$

concentration bound + all-policy realizability  $\not\Rightarrow$  solve batch RL in poly  
(Foster et al. 2021)

$(\frac{1}{\epsilon}, \frac{1}{1-\gamma}, C_{\text{conc}}, \log)$

## Sufficient Conditions (Some results hold for episodic / finite-horizon MDP)

### 1. concentration bound + Bellman completeness / closedness

Given a finite func. class  $\mathcal{F}$ , if  $f \in \mathcal{F}$ ,  $Tf \in \mathcal{F}$

Thm. (Chen & Jiang 2019):  $\sup_{\text{admissible } \pi} \|\frac{v}{\pi}\|_\infty \leq C < \infty$  & completeness, w.p.  $1-\delta$ , the

output policy of FQI after  $K$  iter.,  $\hat{\pi}_K$  satisfies  $\lim_{K \rightarrow \infty} (v^* - v^{\hat{\pi}_K}) \leq \epsilon \cdot V_{\max}$

when  $n = O\left(\frac{C \ln |\mathcal{F}|}{\epsilon^2 (1-\delta)^4}\right)$  discounted  
(mixmax-optimal?)

Why not completeness? ~~stronger than realizability~~

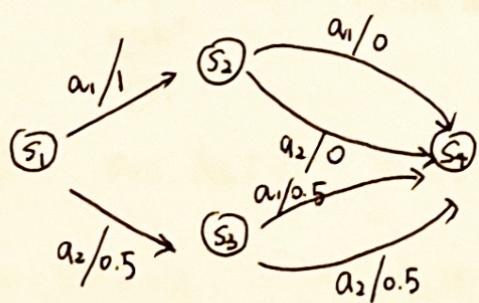
not monotone in  $\mathcal{F}$ , adding more func. may break the condition

(Tkachuk et al. 2024)

### 2. concentration bound + $q^\pi$ -realizability + trajectory data

[Q: how to get a contractive learner  $\hat{\phi}$  /  $\hat{\pi}^\pi$ ] + skipper

target:  $r_h(s, a) + \mathbb{E}_{s' \sim P_h(\cdot | s, a)} [\max_{a'} \phi(s', a')^\top \theta_{h+1}] + \eta_h(s, a) \xrightarrow{\text{noise}} \text{amplified exp.}$



$$q_h^\pi(s', a') \text{ where } \pi = \arg \max_{a'} \phi(s', a')^\top \theta_{h+1}$$

$$\phi(s_1, a_1) = \phi(s_1, a_2) = 1 \quad q(s, a) \text{-realizable}$$

$$\phi(s_2, a_1) = \phi(s_2, a_2) = 0 \quad \text{but } r(s_1, a_1) \& r(s_1, a_2)$$

$$\phi(s_3, a_1) = \phi(s_3, a_2) = 0 \quad \text{is not}$$

$$\sup_{s, a} |r(s, a) - \phi(s, a)^\top \theta| > 0.5/2$$

range(s) =

skip states with approx.  $\sup_a |v^\pi(s) - q^\pi(s, a)|$  small

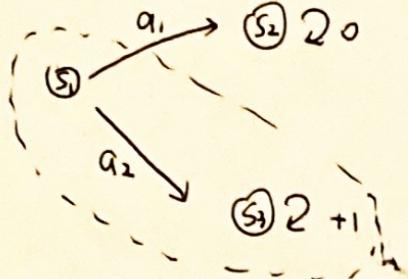
use multi-step rewards  $\Rightarrow \mathbb{E}_{\text{traj}, \text{skip-termination}} [R_h: c_1 + f(s_c)]$

fin. realizable  $\& f: S \rightarrow [0, H]$

stronger than optimality-closedness

$$n = \tilde{\Theta}(C_{\text{conc}}^4 H^7 d^4 / \epsilon^2)$$

3. all-policy-coverage  $\Rightarrow$  single-policy-coverage



$$\phi(s_1) = [0, \frac{1}{1-\gamma}, 0] \quad \theta^* = [1, 1, 0]$$

$$\phi(s_2) = [0, 0, 2] \quad \hat{\theta} = [1, 1, 1]_{\text{initialization}}$$

$$\phi(s_3) = [\frac{1}{1-\gamma}, 0, 0]$$

ensure no over-estimated extrapolation  $\Rightarrow$  instability  $\leftarrow$  pessimism

{ penalized val. func. of under-covered policies  
penalize policies that are far from the one generating data

Vlad's proof  $\forall \pi, \mathbb{E}_{(S_h, A_h) \sim P_h^\pi(\cdot, \cdot)} [\|\phi(S_h, A_h)\|_{X_h^{-1}}] \leq \mathbb{E}_{A_h} [\|\phi\|_{X_h^{-1}}] \cdot C_{\text{conc}} \leq \tilde{O}\left(\frac{C_d}{\sqrt{n}}\right)$

$$X_h \triangleq \sum_{i=1}^{n_h} \phi(S_h^i, A_h^i) \phi(S_h^i, A_h^i)^T + \lambda I$$

$\Rightarrow$  single-policy-concentration  
[finite-horizon]  $\Phi \frac{d_h^\pi(S_h, A_h)}{d_h^u(S_h, A_h)} \leq C \quad \forall h, (S_h, A_h) \text{ for } \pi$

(Jin et al. 2021)  $\mathbb{E}_{h=1}^H [\|\phi(S_h, A_h)\|_{X_h^{-1}}] \leq C(\pi, D)$

(discounted)  $\Phi \max_{f \in \mathcal{F}} \frac{\|f - \gamma^\pi f\|^2 \mathbb{E}_{(s,a) \sim d^\pi} [(f(s,a) - \gamma^\pi f(s,a))^2]}{\mathbb{E}_{(s,a) \sim d^\pi} [(f(s,a) - \gamma^\pi f(s,a))^2]}$

$$\Leftarrow C(\pi, M, \mathcal{F}) \leq \left\| \frac{d^\pi}{M} \right\|_\infty$$

Jin et al. (2021)

single-policy + linear MDP

like UCB  $\hat{W}$  LSVI

$$T_h(s, a) = \beta \|\phi(s, a)\|_{X_h^{-1}} \quad \text{not lin. realizable}$$

$$\hat{Q}_h(s, a) = \min \left\{ \phi(s, a)^T \hat{W} - T_h(s, a), H - h + 1 \right\}^+$$

Xie et al. (2021)  
Zanette (2021)

single-policy +  $\pi$ -Bellman restricted completeness

$$\forall \pi_{h+1} \in \Pi_{h+1}, \forall Q_{h+1} \in \mathcal{F}_{h+1}$$

$$\inf_{Q_h \in \mathcal{F}_h} \|Q_h - \gamma P_h^{\pi_{h+1}} Q_{h+1}\|_\infty = 0$$

Xie:  $\hat{\pi} = \operatorname{argmax}_{\pi \in \Pi} \min_{f \in \mathcal{F}_{\pi, \Sigma}} f(s_0, \pi)$  s.t.  $\mathcal{F}_{\pi, \Sigma} \approx \{f \in \mathcal{F} : \text{empirical Bellman err. under } \pi \leq \varepsilon\}$

Zanette: (Actor Critic) For  $t=1, \dots, T$ ,

$$(\hat{f}_t, \underline{w}_t) = \underset{f \in (\mathbb{R}^d)^H, w \in \mathbb{R}^d}{\operatorname{argmin}} \sum_a \pi_t(a|s_t) \langle \phi(s, a), w \rangle$$

$$\text{s.t. } w_h = \hat{f}_h + \text{LSVI-est. } \hat{W}_h, \| \hat{f}_h \|_{X_h}^2 \leq \lambda_h^2 \text{ & } \| W_h \|_2^2 \leq (\rho_h^w)^2$$

$$\theta_{t+1} = \theta_t + \eta \underline{w}_t$$

$$\pi_h(\theta_{t+1})(a|s) \propto \pi_h(\theta_t)(a|s) e^{\frac{T - \langle \phi(s, a), w \rangle}{\eta} \hat{Q}_h(s, a)}$$

requiring  $\pi$  soft Bellman restricted closedness

Golowich & Moitra (2024) (not peer-reviewed)

lin. Bellman completeness low inherent Bellman err.

$$B_h = \{w \in \mathbb{R}^d : |\langle \phi(s_h, a_h), w \rangle| \leq 1, A(s_h, a)\}$$

$$\text{A1: } \sup_{w \in B_{h+1}} \inf_{w \in B_h} \sup_{\substack{(s, a) \\ \in S_h \times A}} |\langle \phi(s_h, a), w \rangle - \mathbb{E}_{S_{h+1}} [r_h(s_h, a) + \max_{a'} \langle \phi(s_{h+1}, a'), w \rangle]| \leq \frac{\zeta_{BE}}{2}$$

Thm:  $\forall \pi \in \Pi^\delta$  (perturbed linear),  $h \in [H-1]$ ,

$$\sup_{\substack{w \in \mathbb{R}^d \\ a \in A}} \inf_{w_h} |\langle \phi(s_h, a), w_h \rangle - \mathbb{E}_{\substack{S_{h+1} \\ a_{h+1} \sim \pi}} [\langle \phi(s_{h+1}, a_{h+1}), w_{h+1} \rangle]| \leq \tilde{o}( \|w\| d^{3/2} (\sqrt{d} + \frac{1}{\delta}) \zeta_{BE})$$

Def:  $\sigma > 0, h \in [H], \pi_{h, w, \sigma}(s)(a) = \mathbb{P}_{\theta \sim N(w, \sigma^2 I_d)} (a \in \operatorname{argmax} \langle \theta, \phi(s, a) \rangle)$

Thm (informal): A2:  $r_h(s, a) = \langle \phi(s, a), \theta_h \rangle$  determ. & lin. reward

$$\|\phi\| \leq 1, \|\theta_h\| \leq 1, \|w\| \leq B$$

$$\zeta_{BE} \leq C (BH)^2 d^{-3} \quad \text{A3: } C(\pi^*, D)$$

$$\exists \text{ algo. outputs } \hat{\pi}_T, \text{ w.p. } \geq 1 - \delta, \quad V_i^*(s_i) - V_i^{\hat{\pi}_T}(s_i) \leq \tilde{o}(d^{3/2} BH \sqrt{\zeta_{BE}} + \frac{BHd}{\sqrt{n}})(H + C(\pi^*, D)) + \sqrt{\frac{C4B^2 H^2 \sqrt{d}}{T}}$$