

March 25,

Batch RL vs, A

$\forall \pi_{\text{batch}} \in \text{ML}(S, A)$

π -designs

Thm

$\forall \varepsilon, \delta$

$\forall A$

sound

$H = H_{\varepsilon, \delta}$

$\exists M = (S, [A], P, r, r)$

$n(\pi, M, \varepsilon, \delta) \geq c \cdot A^H$

$S \geq H_{\varepsilon, \delta}$

$d = \underline{SA} = HA$
 $= \text{Co}\left(\frac{d}{H}\right)^H$

$Z = (z_1, z_2, \dots, z_n, \dots)$

$z_n \in (S \times A)^n$

} model-based

$\boxed{Z = S \times A}$

$\text{poly}(S, A, H, \frac{1}{\varepsilon^2})$

$\varphi: S \rightarrow \mathbb{R}^d$

$\varphi: S \times A \rightarrow \mathbb{R}^d$

$\rightarrow q^\pi$ -realizability: (M, φ) $\forall \pi: q^\pi \in \mathcal{F}_\varphi$ ✓

q^* -realizability:

$q^* \in \mathcal{F}_\varphi$ X

r^π

r^*

Planning

Batch Learning

~~Batch~~

Adaptively choosing (SA)

Passively choosing T/GA

1. Way - Foster - Kakade $\xrightarrow{\text{Amortila - Jörg Xic -}}$ 2. Andrea Zanette

Theorem: $\forall d > 0, H > 0$ (horizon)

$$P = (P_h)_{h=1}^H - \text{def.}$$

OPE

$\exists (M, \Phi), M = (S, A, P, r); \Phi: S \times A \rightarrow \mathbb{R}^d, \|\Phi\|_2 \leq 1$

$$\text{soES} = r \in [0, 1]$$

$$|S| = O(dH), |A| = 2$$

raw raw
raw
design $\Rightarrow \exists \nu = (\nu_1, \dots, \nu_H), \nu_i \in M_1(S \times A)$

s.t.

- (1) $q_h^\pi \in \mathcal{F}_\Phi = \text{span}(\Phi) \quad \forall h, \forall \pi \in \text{ML}$
- (2) $\mathbb{E}_{P_h} [\Phi \Phi^T] \succeq \frac{1}{d} I \quad \forall h \in [H] \quad (\text{good design})$

(3) $\forall A$ alg. OPE sound

$\forall \pi \in \text{ML}$

$$n(A, M, \Phi, \pi) = \sum \left(\left(\frac{d}{2} \right)^H \right)$$

$\boxed{\pi\text{-design ?}}$

Theorem: $\forall d > 0, H > 0 \quad \exists (M, \Phi), M = (S, A, P, r)$

OPE.

$\Phi: S \times A \rightarrow \mathbb{R}^d, \|\Phi\|_2 \leq 1, r \in [0, 1]$

$$|A| = 2, |S| = O(dH)$$

$\exists \pi_{\text{beh}}, \exists \pi_{\text{trg}} \quad \text{s.t.}$

$$\pi_{\text{beh}} \rightarrow \nu_h^{\text{beh}}, h \in [H]$$

(1) $q_h^{\pi_{\text{trg}}} \in \mathcal{F}_\Phi \quad \forall h$

(2) $\mathbb{E}_{P_h} [\Phi \Phi^T] \succeq \frac{1}{d} I, \forall h$

③ $\forall A$ sound

$$n(A, M, \psi, \pi_{\text{trg}}, \pi_{\text{beh}}) = \sum \left(\left(\frac{d}{2} - 1 \right)^H \right)$$

DPO? \leq ? OPE

Prop:

$$n_{\text{DPO}} \geq n_{\text{OPE}}$$

| \exists -designs
 π -design.

A is "good" at DPO

$\Rightarrow A'(A)$ is "good" at OPE

q^{π} realizability

D : data — MDP M

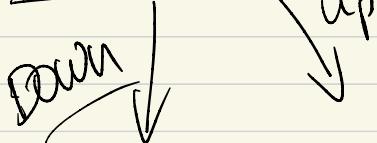
π_{trg} : targ. pol.

s_0 ES: initial

$$V^{\pi_{\text{trg}}}(s_0)$$

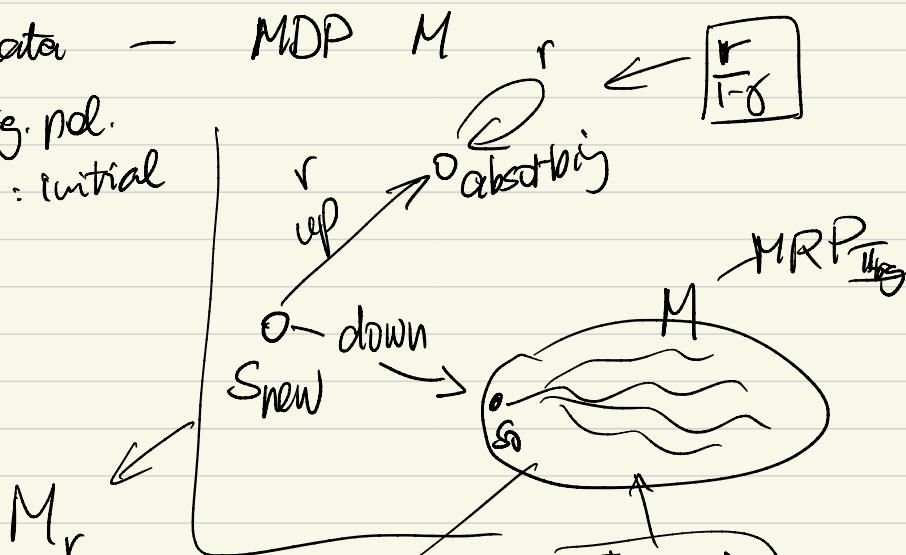
$$r \in [0, 1]$$

$$r := r_0 = \frac{1}{2}$$



$$V^{\pi_{\text{trg}}}(s_0) > \frac{1}{2} \cdot \frac{1}{1-\gamma}$$

$$r_0 = \frac{3}{7}$$



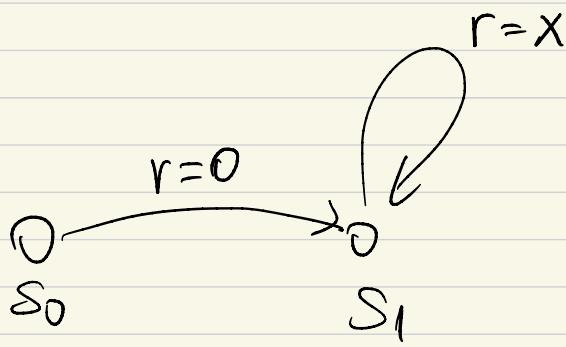
$$(s, a) \\ (s, \pi_{\text{trg}}(s))$$



$$q^{\star} \\ q^{\pi_{\text{trg}}(a)}$$

Discounted unf. horizon

Z-design



$$q^* = v^* = \left[\frac{\gamma x}{1-\gamma} \quad , \quad \frac{x}{1-\gamma} \right] = \underbrace{\begin{bmatrix} \gamma & 1 \end{bmatrix}}_{\Phi} \frac{x}{1-\gamma}$$

$$\varphi(S_0) = \gamma$$

$$\varphi(S_1) = 1$$

Obs : from S_0
only!

$$\theta = \frac{x}{1-\gamma} \leftarrow$$

$$\mathbb{E} [\varphi(S) \varphi(S)]^T = \gamma^2$$

$$P(S=S_0) = 1$$

unidentifiability

π-design

??

Andrea Zanette

$S \sqcup A$ fixed.

π -design

$$g = (g_i)_{i=1}^{\infty}$$

$$\sum_n =$$

$$g_i = (\pi_{11}^{(i)}, \dots, \pi_{k_n}^{(i)}, c_1^{(i)}, \dots, c_{k_n}^{(i)})$$

$$\{(s_1, a_1), \dots, (s_n, a_n)\}$$

$$k_n = n$$

$$s_1^{(i)} = s_1, \dots$$

$$s_n^{(i)} = s_n$$

$$\pi_1^{(i)}(s_1) = a_1$$

:

$$\pi_n^{(i)}(s_n) = a_n$$

$$c_1^{(i)} = \\ = c_n^{(i)} = 1$$



$$Z \subseteq S \times A$$

Observation

$$r|_Z$$

$$P|_Z$$

$$\xrightarrow{S} (S_1^{(i)}, \dots, S_{k_n}^{(i)})$$

$$\text{fix } n \quad \tilde{g} := g_n$$

Z_n = set of readable states

$$\sum_{j=1}^{k_n} c_j^{(i)} \leq n$$

Theorem 1: OPE fixed target.

$\forall r, d \exists S, A$ s.t. $\forall g$ T-design

$\forall V$ alg sound

$\exists (M, \Psi, \Pi_{\text{reg}})$

$M = (S, A, P, r)$

$\|U\|_2 \leq 1, U: S \times A \rightarrow \mathbb{R}^d$

[M, 1]

s.t. $q^{\Pi_{\text{reg}}} \in \mathcal{F}_\Psi$

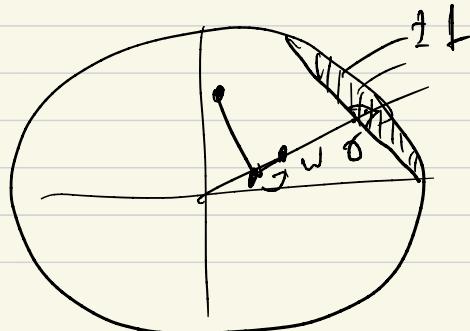
$n(A, M, \Psi, \Pi_{\text{reg}}, g) =$

$$= \mathcal{O}\left(\sqrt{d} \left(\frac{1}{2(1-r)}\right)^d\right)$$

$$\begin{matrix} A^H \\ d^H \end{matrix}$$

$$= \mathcal{O}\left(\sqrt{d} \left(\frac{H}{2}\right)^d\right)$$

$S = A = B_2^d = B = \{x \in \mathbb{R}^d \mid \|x\|_2 \leq 1\}$



$$f(S, a) = a \quad \varrho(S, a) = a$$

Π_{reg}

$$\mu(v^* - v^{\hat{\pi}}) \leq f(n, \nu^*, M_1, \dots)$$

$$\mu v^{\hat{\pi}} \geq \underbrace{\mu v^* - f(n, \nu^*, M_1, \dots)}$$

$$F(n, \nu^*, M_1, \dots)$$