# Lower bounds for OPE w/ linear realizability and uniform coverage.

<u>Recap</u>

- Linear MDP $M = (S, A, H, \{r_h\}_{h \in [H]}, \{P_h\}_{h \in [H]})$

$$r_h(s,a) = \phi(s,a)^T \theta_h, \quad P_h(s'|s,a) = \langle \phi(s,a), \mu_h(s') \rangle$$

- <u>D-optimal design covariance</u> :

$$\Sigma = \underset{(s,a) \sim \rho}{\mathbb{E}} \left[ \phi(s,a) \phi(s,a)^T \right]$$

where $\rho \in \arg\max \quad \log \det \left( \underset{(s,a) \sim U}{\mathbb{E}} \left[ \phi(s,a) \cdot \phi(s,a)^T \right] \right)$

$$U \in \Delta(X \times A)$$

$$|\text{supp}(\rho)| \leq \frac{d(d+1)}{2}$$

property: $\forall x \in \mathbb{R}^d, \quad x^T \Sigma^{-1} x \leq d$

- <u>Offline RL from exploratory datasets</u>

$$H \text{ datasets} \quad D_h = \left\{ \left( s_i, a_i, s_i', r(s_i, a_i) \right) \right\}_{i=1}^{N}$$

Given $i, h$: independent samples $\quad s_i' \sim P_h(\cdot \mid s_i, a_i)$

<u>Assumption</u> "the dataset is exploratory over all dimensions"

$$\frac{1}{N} \sum_{i=1}^{N} \phi(s_i, a_i) \, \phi(s_i, a_i)^{\top} \succcurlyeq \frac{1}{\kappa} \Sigma$$

where $\Sigma$ is the D-optimal design covariance

<u>Alg:</u> let $\hat{\pi} = LSVI\left( \{D_h\}_{h \in [H]}, \phi \right)$

Then:

$$V_1^{\pi^*}(s_1) - V_1^{\hat{\pi}}(s_1) \leq \varepsilon \qquad \text{if } N \geq \frac{poly(H, d, \log(1/\delta), \kappa)}{\varepsilon^2}$$

w.p. $1-\delta$

Today    what are necessary conditions for generalizations?

In particular, we will evaluate realizability in RL
(RL w/ H = 1)

In supervised learning $\overset{\vee}{}$, realizability is sufficient

e.g. PAC w/ 0-1 loss:    $\text{risk}(h_{ERM}) = \Theta\left(\dfrac{d_{VC}}{n}\right)$

w/ realizability

informally, in RL, only realizability is not sufficient for

existence of an sample-efficient algorithm, in the information-theoretic

(minimax) sense.

- <u>Offline policy evaluation (OPE) problem</u>   (perhaps the simplest problem in all RL problems)

  Given $\pi: S \longrightarrow \Delta(A)$ and a feature mapping $\phi: S \times A \longrightarrow \mathbb{R}^d$, the goal:

  output an accurate estimate of $V^\pi$ using collected data sets $\{D_h\}_{h=1}^H$

  using as few samples as possible.

- <u>Realizability assumption</u>   (R)

$$\forall \pi : S \longrightarrow \Delta(A), \; \exists \; \theta_1^\pi, \ldots, \theta_H^\pi \in \mathbb{R}^d :$$

$$Q_h^\pi(s,a) = \phi(s,a)^T \theta_h^\pi \quad \forall (h,s,a)$$

- <u>Data coverage assumption</u>   (strongest possible)   (D)

$$\mathbb{E}_{(s,a) \sim \mu_h} \left[ \phi(s,a) \phi^T(s,a) \right] = \frac{1}{d} I$$

($\mu_h$ satisfy D-optimal design)

- **Theorem**   Assume ($R$). For any algorithm that takes as input both a policy $\pi$ and a feature mapping. $\exists$ MDP satisfying ($D$) s.t:

$$\forall \pi : S \rightarrow \Delta(A), \text{ the algorithm requires } \Omega\left(\left(\frac{d}{L}\right)^{H}\right) \text{ samples}$$

to output the value of $\pi$ up to a constant additive error w/p a.l. 0.9

- **Lemma** (Distinguish Bernoulli random variables)

  - $\alpha \sim \text{Unif}\left(\{\alpha^+, \alpha^-\}\right)$ where

  $$\alpha^- = \frac{1}{2} - \frac{\varepsilon}{2}$$
  $$\alpha^+ = \frac{1}{2} + \frac{\varepsilon}{2}$$

  - $x_1, \ldots, x_n \overset{iid}{\sim} \text{Ber}(\alpha)$

  - $\forall f : \{0,1\}^n \longrightarrow \{\alpha^-, \alpha^+\}$

  $$\Pr\left(f(x_1, \ldots, x_n) \neq \alpha\right) > \underbrace{\frac{1}{4}\left(1 - \sqrt{1 - \exp\left(\frac{-n\varepsilon^2}{1 - \varepsilon^2}\right)}\right)}_{\delta \in (0, 0.25)}$$

  $$n = \frac{1 - \varepsilon^2}{\varepsilon^2} \ln\left(\frac{1}{8\delta(1-\delta)}\right)$$

  If:

  $$0.9 \leq \Pr\left(f(x_1, \ldots, x_n) = \alpha\right) < 1 - \delta$$
  $$\implies \delta \geqslant 0.1 \implies n \geqslant \Omega\left(\frac{1 - \varepsilon^2}{\varepsilon^2}\right)$$

# Hard instances of MDP

Goal:
- construct MDPs that satisfy realizability (R)
- construct offline data that satisfy (D)
- reduction to testing problems

Figure contents (top left diagram):

$\phi(s_h^c, a_1) = e_c$
$\phi(s_h^c, a_2) = e_{c+\hat{d}}$
$\phi(s_h^{\hat{d}+1}, a) = (e_1 + e_2 + \cdots + e_{\hat{d}})/\hat{d}^{1/2}$

Legend:
→ $a_1$
--→ $a_2$

start here

Row 1:
$Q(s, a_1) = r_0 \hat{d}^{(H-1)/2}$
$R(s, a) = 0$

$Q(s^{\hat{d}+1}, a) = r_0 \hat{d}^{H/2}$
$R(s^{\hat{d}+1}, a) = r_0(\hat{d}^{H/2} - \hat{d}^{(H-1)/2})$

Row 2:
$Q(s, a_1) = r_0 \hat{d}^{(H-2)/2}$
$R(s, a) = 0$

$Q(s_2^{\hat{d}+1}, a) = r_0 \hat{d}^{(H-1)/2}$
$R(s_2^{\hat{d}+1}, a) = r_0(\hat{d}^{(H-1)/2} - \hat{d}^{(H-2)/2})$

Row 3:
$Q(s, a_1) = r_0 \hat{d}^{(H-h)/2}$
$R(s, a) = 0$

$Q(s_h^{\hat{d}+1}, a) = r_0 \hat{d}^{(H-h+1)/2}$
$R(s_h^{\hat{d}+1}, a) = r_0(\hat{d}^{(H-h+1)/2} - \hat{d}^{(H-h)/2})$

Row 4:
$Q(s, a_1) = r_0 \hat{d}^{1/2}$
$R(s, a) = 0$

$Q(s_{H-1}^{\hat{d}+1}, a) = r_0 \hat{d}$
$R(s_{H-1}^{\hat{d}+1}, a) = r_0(\hat{d} - \hat{d}^{1/2})$

Row 5:
$Q(s, a) = r_0$
$\mathbb{E}[R(s, a)] = r_0$

$Q(s_H^{\hat{d}+1}, a) = r_0 \hat{d}^{1/2}$
$R(s_H^{\hat{d}+1}, a) = r_0 \hat{d}^{1/2}$

Handwritten notes (right side):

$$Q_h^\Pi(s_h^c, a_1) = r_0 \, \hat{d}^{(H-h)/2}$$

$$Q_h^\Pi(s_h^c, a_2) = r_0 \, \hat{d}^{(H-h)/2}$$

$$Q_h^\Pi(s_h^{\hat{d}+1}, a) = r_0 \, \hat{d}^{(H-h+1)/2}$$

$$Q_h^\Pi(s, a) = \phi(s, a)^T w_h^\Pi$$

$$R(s_{H,1}^c, a) = \begin{cases} 1 & \text{wp} \quad \dfrac{(1+r_0)}{2} \\ -1 & \text{wp} \quad \dfrac{1-r_0}{2} \end{cases}$$

Handwritten notes (bottom):

$$Q_h^\Pi(s, a) = \phi(s, a)^T \begin{bmatrix} r_0 \hat{d}^{(H-h)/2} \\ \vdots \\ (H-h)/2 \\ r_0 \hat{d} \\ \vdots \end{bmatrix} \Big\} \hat{d}$$

- **Dataset:** $\mu_h$ is uniform over $\{(s_h^c, a_1), (s_h^c, a_2)\}_{c \in [\hat{d}]}$

$$\mathbb{E}_{(s,a) \sim \mu_h} \left[ \phi(s,a) \phi(s,a)^T \right] = \frac{1}{d} I$$

- **Reduction to testing:**

  - the policy value to estimate: $V_1^\pi(s_1^{\hat{d}+1}) = r_0 \, 2^{H/2}$

  - consider 2 instances: (I) $r_0 = 0 \longrightarrow V_1^\pi(s_1^{\hat{d}+1}) = 0$

    (II) $r_0 = 2^{-H/2} \longrightarrow V_1^\pi(s_1^{\hat{d}+1}) = 1$

  - if the algorithm wants to output an estimate that is correct up to $0.5$ error, then it must need to distinguish two problem instances (I) and (II)

**Consider any algorithm:**

$$\text{Alg}: \left( \{P_h\}_{h \in [H]}, \phi \right) \longrightarrow \mathbb{R}$$

- Given datasets $\{D_h\}_{h \in [H]}$ (and $\Phi$), the algorithm need to identify which of the two instances (I) and (II) that the datasets come from

- Note that for both (I) and (II):
  - data distribution, transition kernels are the same
  - rewards are zero everywhere except in the last layer $H$

- Thus, to distinguish (I) and (II), the algorithm need to distinguish the reward distribution:

$$r = \begin{cases} 1 & \text{wp } \frac{1}{2} \\ -1 & \text{wp } \frac{1}{2} \end{cases}$$

and

$$r = \begin{cases} 1 & \text{wp } \dfrac{1 + \bar{d}^{-H/2}}{2} \\ -1 & \text{wp } \dfrac{1 - \bar{d}^{-H/2}}{2} \end{cases}$$

$$\Rightarrow \quad n = \Omega\left( \hat{d}^{H} \right) = \Omega\left( \left(\frac{d}{2}\right)^{H} \right)$$