

$\Delta : \mathcal{X} \times \Theta \rightarrow \mathbb{R}$ - discriminant fn.

Def 1: $x_{1:n} \in \mathcal{X}^n$ is $\varepsilon_1/\varepsilon_2/\ell^\infty$ -cluster sequence if
 $\forall i \in [n] \exists \theta \in \Theta$ s.t. $|\Delta(x_i, \theta)| > \varepsilon_1$ & $\max_{j \leq i-1} |\Delta(x_j, \theta)| \leq \varepsilon_2$

Set of these: $S_\infty(\varepsilon_1, \varepsilon_2)$.

Def 2: $x_{1:n} \in \mathcal{X}^n$ is $\varepsilon_1/\varepsilon_2/\ell^2$ -cluster seq. if
 $\forall i \in [n] \exists \theta \in \Theta$ s.t. $|\Delta(x_i, \theta)| > \varepsilon_1$ & $\sum_{j=1}^{i-1} \Delta^2(x_j, \theta) \leq \varepsilon_2^2$

Set of these: $S_2(\varepsilon_1, \varepsilon_2)$.

① S_∞, S_2 are decreasing in ε_1 , increasing in ε_2 .

② $S_\infty(\varepsilon_1, \varepsilon_2) \cap \mathcal{X}^n \subseteq S_2(\varepsilon_1, \sqrt{n} \varepsilon_2)$, $\forall \varepsilon_1, \varepsilon_2, n$. [$\theta_1, \dots, \theta_{n-1}$ witness for ℓ^∞
 $\Rightarrow \sum_{j=1}^{i-1} \Delta^2(x_j, \theta_i) \leq (i) \varepsilon_2^2 \leq (n) \varepsilon_2^2$.

Claim: $\max \{n \mid S_\infty(\varepsilon, \frac{\varepsilon}{\sqrt{d}}) \cap \mathcal{X}^n \neq \emptyset\} \leq \max \{n \mid S_2(\varepsilon, \varepsilon) \cap \mathcal{X}^n \neq \emptyset\} =: d$.

Proof: It suffices if $S_\infty(\varepsilon, \frac{\varepsilon}{\sqrt{d}}) \cap \mathcal{X}^{d+1} = \emptyset$. And for this: $S_\infty(\varepsilon, \frac{\varepsilon}{\sqrt{d}}) \cap \mathcal{X}^{d+1} \subseteq S_2(\varepsilon, \varepsilon) \cap \mathcal{X}^{d+1} = \emptyset$ // Q.e.d.

We can directly bound $\max \{n \mid S_\infty(\varepsilon_1, \varepsilon_2) \cap X^n \neq \emptyset\}$

when $\Delta(x, \theta) = x^T \theta$, $X = B_2^d(\sigma)$, $\Theta = B_2^d(S)$.

How? Usual proof: Let $x_{1:n} \in S_\infty(\varepsilon_1, \varepsilon_2)$. Take $i \in [n]$.

$$\Rightarrow \varepsilon_1 < \max \{ |x_i^T \theta| : \max_{1 \leq j \leq i-1} |x_j^T \theta| \leq \varepsilon_2, \|\theta\|_2 \leq S \}$$

$$\leq \max \{ |x_i^T \theta| : \sum_{j=1}^{i-1} |x_j^T \theta|^2 \leq (i-1)\varepsilon_2^2, \|\theta\|_2^2 \leq S^2 \}$$

$$= \max \{ |x_i^T \theta| : \theta^T (X_{i-1} + \lambda I) \theta \leq i \varepsilon_2^2 \}$$

$$= \sqrt{i} \varepsilon_2 \|x_i\|_{V_{i-1}^{-1}} \Rightarrow \|x_i\|_{V_{i-1}^{-1}} > \frac{\varepsilon_1}{\sqrt{i} \varepsilon_2}$$

$$\left(\frac{\lambda d + i \sigma^2}{d} \right)^d = \left(\frac{\text{tr} V_i}{d} \right)^d \geq \det V_i = \lambda^d \prod_{j=1}^i \left(1 + \|x_j\|_{V_{j-1}^{-1}} \right) = \lambda^d \prod_{j=1}^i \left(1 + \frac{\varepsilon_1}{\sqrt{j} \varepsilon_2} \right)$$

$$\Rightarrow d \log \left(1 + \frac{i \sigma^2}{d \lambda} \right) \geq \sum_{j=1}^i \log \left(1 + \frac{\varepsilon_1}{\sqrt{j} \varepsilon_2} \right) \approx \frac{\varepsilon_1}{\varepsilon_2} \sum_{j=1}^i \frac{1}{\sqrt{j}} \approx \frac{\varepsilon_1}{\varepsilon_2} \sqrt{i}$$

$$\Rightarrow i_{\max} = \mathcal{O}(d^2 \log^2(\dots))$$

$$X_i = \sum_{j=1}^i x_j x_j^T$$

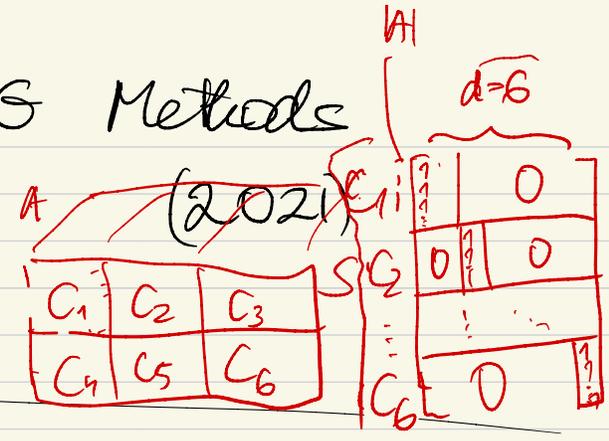
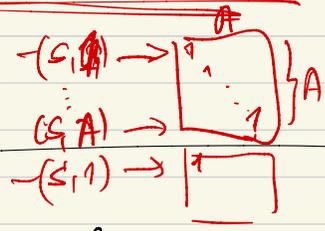
$$\lambda = \varepsilon_2^2 / S^2$$

$$V_i = X_i + \lambda I$$

solve for largest i

Approximate Benefits of PG Methods with Aggregated States

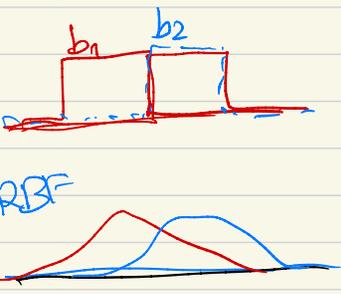
Dan Russo



Lower bound for API

$$\tilde{\epsilon}_{\text{apx}} = \sup_{\pi \in \Pi_{\Phi}} \inf_{\theta} \|\Phi \theta - q^{\pi}\|_{\infty}$$

Theorem: $\forall \gamma \in [0, 1), \forall \epsilon_{\text{apx}} > 0$ \exists MDP $M = (S, A, P, r, \gamma)$
 $\exists \Phi$
 $\exists \pi_1$ policy of M ; initial API policy

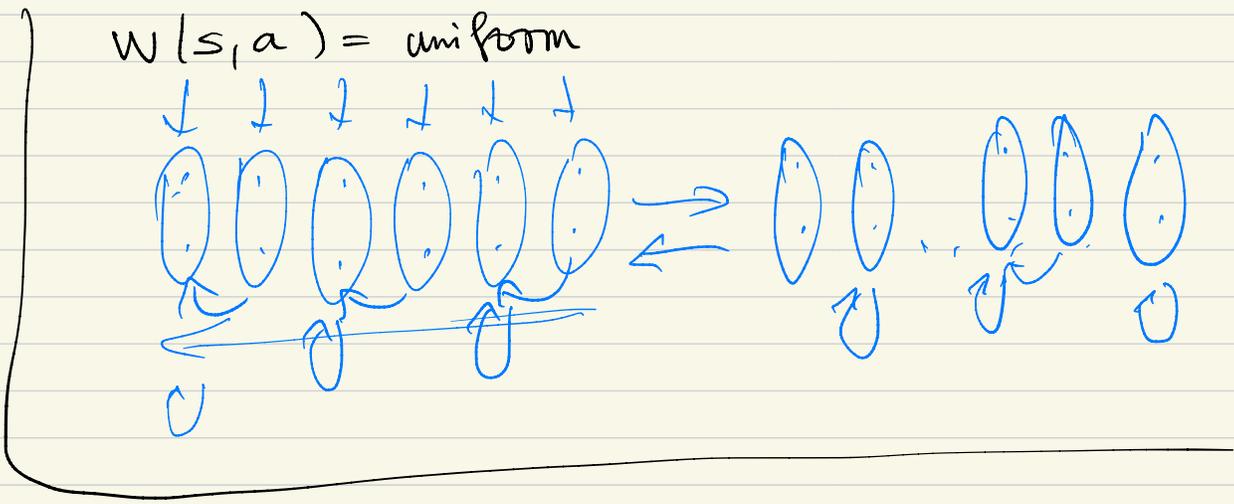


$\exists g \in \mathcal{M}_1(S)$

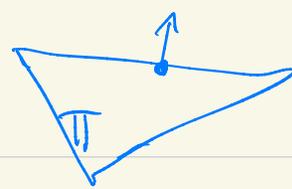
$\rightsquigarrow \pi_1, \pi_2, \pi_3, \dots$

s.t. $\inf_{t \geq 1} v^*(g) - v^{\pi_t}(g) \geq \frac{c \cdot \tilde{\epsilon}_{\text{apx}}}{(1-\gamma)^2}$

Policy eval: $\hat{q}_t = \underset{q \in \mathcal{F}_{\Phi}}{\text{argmin}} \sum_{\substack{S \in S \\ a \in A}} w(s,a) (q(s,a) - q^{\pi_t}(s,a))^2$



$$J(\pi) \doteq V^\pi(g) \rightarrow \max$$



Thm: π^∞ stat. point of $J(\pi)$

$$\frac{1}{1-\gamma} \left(V^*(g) - V^\pi(g) \right) \leq \frac{\gamma \epsilon_{\text{apx}}}{1-\gamma}$$

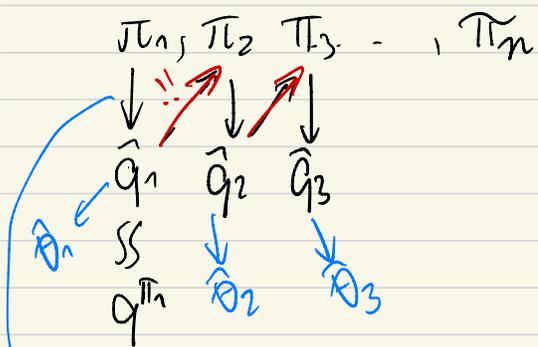


$$(1-\gamma)(V^*(g) - V^\pi(g)) \leq \gamma \epsilon_{\text{apx}}$$

$$\tilde{V}^*(g) - \tilde{V}^\pi(g) \leq \gamma \epsilon_{\text{apx}}$$

POLITEX

Policy Iteration with Expert Advice



$$\pi_k(a|s) \propto \exp\left(\eta \sum_{j=1}^{k-1} \hat{q}_j(s, a)\right) = E_k(sa)$$

$$\pi_k(a|s) = \frac{E_k(sa)}{\sum_{a'} E_k(sa')}$$

LSPE

G-optimal design
m rollouts, H length,

$$\bar{\theta}_{k-1} = \sum_{j=1}^{k-1} \hat{\theta}_j$$

$$\hat{q}_k = \Phi \hat{\theta}_k$$

$$\sum \hat{q}_j = \Phi \sum \hat{\theta}_j$$

Why does it work?

$$\pi_1, \dots, \pi_n \quad \frac{1}{n}(\pi_1 + \dots + \pi_n)$$

$$K \sim \text{Unif}([n])$$

$$A = \pi_K(s_0)$$

Politer in plain

π : policy induced by Politer.

$$R = \sum_{t=0}^{\infty} \gamma^t r_{A_t}(s_t)$$

$$\underline{v^\pi(s)} = \mathbb{E}_s^\pi [R] \stackrel{\downarrow}{=} \frac{1}{n} \sum_{k=1}^n \mathbb{E}_s^{\pi_k} [R] = \frac{1}{n} \sum_{k=1}^n \underline{v^{\pi_k}(s)}$$

$$P_s^\pi = \frac{1}{n} \sum_{k=1}^n P_s^{\pi_k}$$

$$\frac{1}{n} \sum_{k=1}^n v^{\pi^*} - v^{\pi_k}$$

$$= \frac{1}{n} (\mathbb{I} - \gamma P_{\pi^*})^{-1} \sum_{k=1}^n \left[T_{\pi^*} v^{\pi_k} - v^{\pi_k} \right]$$

$$v^{\pi^*} - v^{\pi_k} = (\mathbb{I} - \gamma P_{\pi^*})^{-1} [T_{\pi^*} v^{\pi_k} - v^{\pi_k}]$$

$$(M_\pi q)(s) = \sum_a \pi(a|s) q(s,a)$$

$$T_{\pi^*} v^{\pi_k} = r_{\pi^*} + \gamma P_{\pi^*} v^{\pi_k}$$

$$r_\pi = M_\pi r$$

$$P_\pi = M_\pi P$$

$$= M_{\pi^*} (r + \gamma P v^{\pi_k})$$

$$= M_{\pi^*} q^{\pi_k}$$

$$v^{\pi_k} = M_{\pi_k} q^{\pi_k}$$

$\frac{\text{Exp } \gamma^t}{1-\gamma}$

$\text{Exp } \gamma^t + \dots$

$$= \frac{1}{n} (\mathbb{I} - \gamma P_{\pi^*})^{-1} \sum_{k=1}^n M_{\pi^*} q^{\pi_k} - M_{\pi_k} q^{\pi_k}$$

$$= \frac{1}{n} (\mathbb{I} - \gamma P_{\pi^*})^{-1} \sum_{k=1}^n \left(M_{\pi^*} \hat{q}_k - M_{\pi_k} \hat{q}_k + \frac{1}{n} (\mathbb{I} - \gamma P_{\pi^*})^{-1} \sum_{l=1}^n (M_{\pi^*} - M_{\pi_l}) \times (q^{\pi_l} - \hat{q}_l) \right)$$

I. \rightarrow (?)

II.