

Function approx.

$$\Phi = \begin{bmatrix} \varphi_1(\mathbf{x}) \\ \vdots \\ \varphi_d(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^d$$

$$(B2) \quad \text{iff} \quad q^{\pi} \in \mathcal{F} - \text{span}(\Phi)$$

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

$$\exists \theta \in \mathbb{R}^d : q^{\pi} = \Phi \theta$$

$$[\forall s, a : \Psi(s, a)^T \theta = q^{\pi}(s, a)]$$

query + runtime  
generative model  
indep. #S

P.I.:  $\partial \Psi \theta_k = 0, 1, 2, \dots$  : step greedy w.r.t.  $q^{\pi}_{\theta_k}$

$$\rightarrow q^{\pi_{\theta_k}} = \Phi \theta_k, \quad \theta_k \in \mathbb{R}^d$$

$d \ll \#S$

extrapolation

Linear regression :  $\|\Phi \hat{\theta}_k - q^{\pi_k}\|_{\infty}$  small  
 $\uparrow$   $S \times A \gg d$

KW-Theorem

$$\exists C \subseteq \mathcal{Z} = S \times A$$

$$\exists g: C \rightarrow [0, 1] \quad \sum_{z' \in C} g(z') = 1$$

$$1) |C| \leq \frac{d(d+1)}{2}$$

$$2) \sup_{z \in \mathcal{Z}} \|\Psi(z)\|_{G_g^{-1}} \leq \sqrt{d}$$

error blowup  
factor  
of least-squares.

$$G_g = \sum_{z \in C} g(z) \Psi(z) \Psi(z)^T$$

$$\|u\|_p = (\mathbf{u}^T P \mathbf{u})^{1/2}$$

$\sqrt{d}$  unimprovable

$S \times A$  finite

$\Psi(S \times A) \subseteq \mathbb{R}^d$

compact

$\forall \theta, \varepsilon: \mathbb{Z} \rightarrow \mathbb{R}$

$$\text{Cor: } \hat{\theta} = G_S^{-1} \sum_{z' \in C} g(z') (\varphi(z)^T \theta + \varepsilon(z)) \varphi(z)$$

measured

$$\text{Then } \max_{z \in Z} |\varphi(z)^T \hat{\theta} - \varphi(z)^T \theta| \leq \left( \max_{z' \in C} |\varepsilon(z')| \right) \sqrt{d}$$

$$\text{Proof: } \hat{\theta} = \theta + G_S^{-1} \sum_{z' \in C} g(z') \varepsilon(z') \varphi(z')$$

$z \in Z$  fixed

$$|\varphi(z)^T \hat{\theta} - \varphi(z)^T \theta| \leq \sum_{z' \in C} |\varepsilon(z')| |g(z')| |\varphi(z)^T G_S^{-1} \varphi(z')|$$

$$\leq \underbrace{\left( \max_{z' \in C} |\varepsilon(z')| \right)}_{\varepsilon_{\text{apx}}} \left( \sum_{z' \in C} |g(z')| |\varphi(z)^T G_S^{-1} \varphi(z')| \right)^{1/2}$$

$$\left( \int f d\mu \right)^2 \leq \varepsilon_{\text{apx}} \left( \sum_{z' \in C} |g(z')| \underbrace{|\varphi(z)^T G_S^{-1} \varphi(z')|^2}_{\varphi(z)^T G_S^{-1} \varphi(z')} \right)^{1/2}$$

$$\leq \int f^2 d\mu = \varepsilon_{\text{apx}} \left( \sum_{z' \in C} |g(z')| \underbrace{\varphi(z)^T G_S^{-1} \varphi(z')}_{\varphi(z)^T G_S^{-1} \varphi(z')} \varphi(z)^T G_S^{-1} \varphi(z) \right)^{1/2}$$

$$\leq \int \varphi(f) d\mu = \varepsilon_{\text{apx}} \left( \varphi(z)^T G_S^{-1} \underbrace{\left( \sum_{z' \in C} g(z') \varphi(z)^T G_S^{-1} \varphi(z') \varphi(z)^T G_S^{-1} \varphi(z') \right)}_{G_S} \varphi(z) \right)^{1/2}$$

$$= \mathbb{E}_{\text{apx}} \left( \underbrace{\varphi(z)^T G_S^{-1} \varphi(z)}_{\| \varphi(z) \|_{G_S^{-1}}^2} \right)^{1/2}$$

$$= \mathbb{E}_{\text{apx}} \cdot \underbrace{\| \varphi(z) \|_{G_S^{-1}}}_{\leq \sqrt{d}} \quad // \text{Qu.-e.d.}$$

P.E. (Policy Evaluation)

$$\pi \rightarrow \underline{q^\pi} = \underline{\Phi \theta} + \underline{\varepsilon_\pi}$$

$$\boxed{\hat{\theta} = G_S^{-1} \sum_{z \in C} g(z) \hat{R}_m(z) \varphi(z)}$$
LS

$$\hat{R}_m(s, a) = \frac{1}{m} \sum_{j=1}^m \left[ \sum_{t=0}^H \gamma^t r_{A_t^{(j)}}(S_t^{(j)}) \right]$$

$\underbrace{\text{truncated return}}_{j=1, \dots, m}$



$$S_0^{(j)} = s$$

$$A_0^{(j)} = a \quad t \geq 1$$

$$S_t^{(j)} \sim P_{A_{t-1}^{(j)}}(S_{t-1}^{(j)})$$

$$A_t^{(j)} \sim \pi(\cdot | S_t^{(j)})$$

$$\textcolor{red}{\cancel{\theta}} |q^\pi(z) - \varphi(z)^\top \hat{\theta}| = |\varphi(z)^\top \theta + \varepsilon_\pi(z) - \varphi(z)^\top \hat{\theta}|$$

$$\leq |\varphi(z)^\top \theta - \varphi(z)^\top \hat{\theta}| + \max_{z' \in \mathcal{Z}} |\varepsilon_\pi(z')| \quad \textcolor{red}{\cancel{\varepsilon_\pi^*} \in \mathbb{R}}$$

$$\hat{R}_m(z) = \varphi(z)^\top \theta + \underbrace{\hat{R}_m(z) - \varphi(z)^\top \theta}_{\varepsilon(z)}$$

Using Corollary:

$$|\varphi(z)^\top \theta - \varphi(z)^\top \hat{\theta}| \leq \max_{z \in \mathcal{C}} |\varepsilon(z)| / \sqrt{\alpha}$$

sample error

$$\varepsilon(z) = \underbrace{\hat{R}_m(z)}_{\text{I}} - (\mathbf{T}_\pi^\# \mathbf{O})(z) + (\mathbf{T}_\pi^\# \mathbf{O})(z) - \underbrace{q^\pi(z) - \varepsilon_\pi(z)}_{\text{II} + \text{III}}$$

$$\left[ \mathbb{E}[ \varepsilon] = (\mathbf{T}_\pi^\# \mathbf{O})(z) \right] \leq \gamma^\# / 1 - \gamma$$

$$q^\pi(z) = \varphi(z)^\top \theta + \varepsilon_\pi(z)$$

$$\text{Hoeffding's} \leq |\hat{R}_m(z) - (\mathbf{T}_\pi^\# \mathbf{O})(z)|$$

Hoeffding: wp 1 - δ

$$\leq \frac{1}{1-\gamma} \sqrt{\frac{\log(2/\delta)}{2m}}$$

+ Union bound  $\geq \mathcal{C}$

wp 1- $\delta$

$\forall z \in \mathcal{C}$

$\delta \rightarrow \delta/Cl$

$$|\hat{P}_m(z) - (\Pi^H \theta)(z)| \leq \frac{1}{1-\gamma} \sqrt{\frac{\log(\frac{2|C|}{\delta})}{2m}}$$

wp 1- $\delta$ :  $\forall z$ :

$$|\varepsilon(z)| \leq \frac{1}{1-\gamma} \left( \dots + \frac{\gamma^H}{1-\gamma} + \varepsilon_{\Pi}^* \right)$$

Lemma:

wp 1- $\delta$

$$\max_{z \in \mathcal{Z}} |q^T(z) - \ell(z)^T \theta| \leq \underline{\varepsilon_{\Pi}^*} + \overline{\varepsilon_{\Pi}^*} \left( \varepsilon_{\Pi}^* + \frac{\gamma^H}{1-\gamma} + \frac{1}{1-\gamma} \sqrt{\frac{2|C|}{2m}} \right)$$

$m \rightarrow \infty$   $\theta \rightarrow$

Notice

$$|C| = O(d^2)$$

$$\log |C| = O(\log d)$$

note  
 depends  
 on #S

$$\varepsilon_{\Pi}^* (1+3\sqrt{d})$$

$$\frac{\gamma^H}{1-\gamma} \leq \varepsilon_{\text{approx}}$$

$$H = H_{\gamma, \varepsilon_{\text{approx}}}$$

$$\frac{1}{1-\gamma} \sqrt{\frac{\log(\frac{2|C|}{\delta})}{2m}} \leq \varepsilon_{\text{approx}}$$

$d^2$

$$m \geq \left( \frac{1}{(1-\gamma)\varepsilon_{\text{approx}}} \right)^2 \log \left( \frac{2|C|}{\delta} \right)$$

$$\pi_0 \quad \cancel{q^{\pi_0}} \quad \theta_0 \quad q_0 = \Phi \theta_0 = q^{\pi_0} + \underline{\varepsilon}_0$$

$\pi_1$  greedy w.r.t.  $q_0$  (not  $q^{\pi_0}$ !)

⋮

$$\pi_k - \text{---} \quad q_k = \Phi \theta_k = q^{\pi_k} + \underline{\varepsilon}_k$$

$$V^{\pi_k} \geq V^* - \text{??}_1$$

& large enough

$$\max_{1 \leq i \leq k} \|\varepsilon_i\|_\infty \leq \delta$$

Progress Lemma A.P.I.  $\pi$  ML policy

$$\geq T_{\pi^*} q \quad | \quad q = q^\pi + \varepsilon \quad \varepsilon: S \times A \rightarrow \mathbb{R}$$

$$\underline{\pi'} : \boxed{T_{\pi'} q = T q} \quad (\pi' \text{ greedy wrt. } q)$$

$$\boxed{M_{\pi'} q = M q : \sum_{a \in A} \pi'(a|s) q(s, a) = \max_a q(s, a)}$$

$$\boxed{\|q^* - q^{\pi'}\|_\infty \leq \gamma \|q^* - q^\pi\|_\infty + \frac{2\|\varepsilon\|_\infty}{1-\gamma}}$$

Proof:  $\pi^*$  opt. ML policy  $T_{\pi^*} q^* = q^*$

$$q^* - \boxed{q^{\pi'}} = \underbrace{T_{\pi^*} q^* - T_{\pi^*} q^{\pi'}}_{=0} + \underbrace{T_{\pi^*} q^{\pi'} - T_{\pi'} q^{\pi'}}_{\leq 0} + \underbrace{T_{\pi'} q^{\pi'} - T_{\pi'} q}_{\leq 0}$$

$$+ \underbrace{T_{\pi'} q}_{P} - \underbrace{T_{\pi'} q^{\pi}}_{P} + \underbrace{T_{\pi'} q^{\pi}}_{P} - \underbrace{T_{\pi'} q^{\pi^*}}_{P}$$

$$T_{\pi'} q = r + \gamma \underbrace{P M_{\pi'} q}_{P}$$

$$\leq \sigma P_{\bar{\pi}^*}(q^* - q^\pi) + \sigma P_{\bar{\pi}^*}(q^\pi - q)$$

$$+ \sigma P_{\bar{\pi}^*}(q - q^\pi) + \sigma P_{\bar{\pi}^*}(q^\pi - q^{\pi'})$$

$$= \sigma P_{\bar{\pi}^*}(q^* - q^\pi) + \sigma [(\underbrace{P_{\bar{\pi}^*} - P_{\bar{\pi}^*}}_{\cancel{\text{I}}}) \varepsilon + P_{\bar{\pi}^*}(\underbrace{q^\pi - q}_{\cancel{\text{II}}})]$$

$$q^\pi - q^{\pi'} = (I - \sigma P_{\bar{\pi}^*})^{-1} \left[ \underbrace{(I - \sigma P_{\bar{\pi}^*}) q^\pi - r}_{\substack{\text{I} \\ (I - \sigma P_{\bar{\pi}^*})^{-1} r}} \right] \quad \begin{array}{l} q^\pi = q - \varepsilon \\ \sigma (P_{\bar{\pi}^*} - P_{\bar{\pi}^*}) \varepsilon \end{array}$$

$$\underline{I.} = (I - \sigma P_{\bar{\pi}^*})(q - \varepsilon) - r$$

$$= q - \underbrace{(\sigma P_{\bar{\pi}^*} q + r)}_{(I - \sigma P_{\bar{\pi}^*}) \varepsilon} - (I - \sigma P_{\bar{\pi}^*}) \varepsilon$$

$$= q - \underline{T_{\bar{\pi}^*} q} - (I - \sigma P_{\bar{\pi}^*}) \varepsilon$$

$$q = q^\pi + \varepsilon / T_{\bar{\pi}^*}$$

$$\leq (I - \sigma P_{\bar{\pi}^*}) \varepsilon - (I - \sigma P_{\bar{\pi}^*}) \varepsilon$$

$$T_{\bar{\pi}^*} q = \underline{T_{\bar{\pi}^*} q} + \sigma P_{\bar{\pi}^*} \varepsilon$$

$$= \sigma (P_{\bar{\pi}^*} - P_{\bar{\pi}^*}) \varepsilon.$$

$$q - T_{\bar{\pi}^*} q \leq (I - \sigma P_{\bar{\pi}^*}) \varepsilon$$