Lecture 8 Planning with global access and uniform realizability

### Administrivia: Projects

- What makes a good project?
- Choose a topic
- If you choose a paper, improve on it
  - How does it fit into the big picture?
  - What is left? Can you add anything?
  - Can you make the proof nicer?
- Move on if the paper is not "good"
- Start early
- Ask for help

#### **Global access**

Can get all the features at all states, can preprocess it

#### $\underline{\zeta}$ –uniform action-value realizability

$$\sup_{\pi \in \text{DETML}} \inf_{\theta} \|q^{\pi} - \Phi\theta\|_{\infty} \leq \zeta$$

#### <u>Note</u>

For finite MDPs,  $\forall \pi \in ML$ ,  $\exists m > 0, (\pi_i)_{i \in [m]}$  st.  $q^{\pi} \in \Sigma_i \alpha_i q^{\pi_i}$  with some  $\alpha_i \ge 0, \sum_i \alpha_i = 1$ (Dadashi et al.) least-squares policy evaluation

#### **Least-squares Policy Evaluation**

1. Rollouts from a set C of wellchosen state-action pairs

- 2. Average over those
- 3. Least-squares fit



 $(\mathbb{E}\big[\hat{R}_m(z)\big] = q^\pi(z))$ 

$$\hat{ heta} = rg\min_{ heta \in \mathbb{R}^d} \sum_{z \in \mathcal{C}} arrho(z) \Big( \langle heta, arphi(z) 
angle - \hat{R}_m(z) \Big)^2$$

$$\hat{R}_m(z) = rac{1}{m}\sum_{j=1}^m\sum_{t=0}^{H^{(j)}-1}r_{A_t^{(j)}}(S_t^{(j)})$$

 $\mathbb{P}(H^{(j)} \geq t+1) = \gamma^t$ 

First head with  $X_1, X_2, ...$ :  $\mathbb{P}(X_i = \text{Head}) = 1 - \gamma$ 

$$\hat{ heta} = rg\min_{ heta \in \mathbb{R}^d} \sum_{z \in \mathcal{C}} arrho(z) \Big( \langle heta, arphi(z) 
angle - \hat{R}_m(z) \Big)^2$$

Alternative: Choose *H* large enough and let

$$\hat{R}_m(z) = rac{1}{m}\sum_{j=1}^m\sum_{t=0}^{H-1}\gamma^t r_{A_t^{(j)}}(S_t^{(j)}).$$

Lemma will be for this choice.

Homework: Think about the pros and cons of switching to the trajectories with random length

# $G_{g} = \frac{2g(z)}{zec} \left( \frac{2}{z} \right)^{T}$

**Lemma (extrapolation error control in least-squares):** Fix any  $\theta \in \mathbb{R}^d$ ,  $\varepsilon : \mathbb{Z} \to \mathbb{R}$ ,  $\mathcal{C} \subset \mathbb{Z}$  and  $\varrho \in \Delta_1(\mathcal{C})$  such that the moment matrix  $G_{\varrho}$  is nonsingular. Define

$$\hat{ heta} = G_arrho^{-1} \sum_{z' \in C} arrho(z') \Big( arphi(z')^ op heta + arepsilon(z') \Big) arphi(z') \,.$$

Then, for any  $z \in \mathcal{Z}$  we have

 $\left| \varphi(z)^{\top} \hat{\theta} - \varphi(z)^{\top} \theta \right| \leq \left\| \varphi(z) \right\|_{C_{e}^{-1}} \max_{z' \in C} |\varepsilon(z')|.$ 

#### optimal design

**Theorem (Kiefer-Wolfowitz):** Let  $\mathcal{Z}$  be finite. Let  $\varphi : \mathcal{Z} \to \mathbb{R}^d$  be such that the underlying feature matrix  $\Phi$  is rank d. There exists a set  $\mathcal{C} \subseteq \mathcal{Z}$  and a distribution  $\varrho : C \to [0, 1]$  over this set, i.e.  $\sum_{z' \in \mathcal{C}} \varrho(z') = 1$ , such that

- 1  $|\mathcal{C}| \leq d(d+1)/2;$
- $^{2} \ \, \sup_{z\in\mathcal{Z}}\|\varphi(z)\|_{G^{-1}_{\varrho}}\leq \sqrt{d};$
- <sup>3</sup> In the previous line, the inequality is achieved with equality and the value of  $\sqrt{d}$  is best possible under all possible choices of C and  $\rho$ .

#### extrapolation error control with LS

**Corollary (extrapolation error control in least-squares via optimal design):** Fix any  $\varphi : \mathbb{Z} \to \mathbb{R}^d$  full rank. Then, there exists a set  $\mathcal{C} \subset \mathbb{Z}$  with at most d(d+1)/2 elements and a weighting function  $\varrho \in \Delta_1(\mathcal{C})$  such that for any  $\theta \in \mathbb{R}^d$  and any  $\varepsilon : \mathcal{C} \to \mathbb{R}$ ,

$$\max_{z\in\mathcal{Z}} \left| arphi(z)^ op \hat{ heta} - arphi(z)^ op heta 
ight| \leq \sqrt{d} \, \max_{z'\in C} \left| arepsilon(z') 
ight|.$$

where  $\hat{\theta}$  is given by

$$\hat{ heta} = G_arrho^{-1} \sum_{z' \in C} arrho(z') \Big( arphi(z')^ op heta + arepsilon(z') \Big) arphi(z') \,.$$

C and  $\varrho$  are chosen independently of  $\theta$  and  $\epsilon$ ! **Lemma (LSPE-***G* **extrapolation error control):** Fix any full-rank feature-map  $\varphi : \mathbb{Z} \to \mathbb{R}^d$  and take the set  $\mathcal{C} \subset \mathbb{Z}$  and the weighting function  $\varrho \in \Delta_1(\mathcal{C})$  as in the Kiefer-Wolfowitz theorem. Fix an arbitrary policy  $\pi$  and let  $\theta$  and  $\varepsilon_{\pi}$  such that  $q^{\pi} = \Phi \theta + \varepsilon_{\pi}$  and assume that immediate rewards belong to the interval [0, 1]. Let  $\hat{\theta}$  be as in Eq. (6). Then, for any  $0 \leq \delta \leq 1$ , with probability  $1 - \delta$ ,

$$\left\|q^{\pi} - \Phi\hat{\theta}\right\|_{\infty} \leq \|\varepsilon_{\pi}\|_{\infty} (1 + \sqrt{d}) + \sqrt{d} \left(\frac{\gamma^{H}}{1 - \gamma} + \frac{1}{1 - \gamma} \sqrt{\frac{\log(2|C|/\delta)}{2m}}\right).$$
(7)

KW:  $|\mathcal{C}| \leq d(d+1)/2$ , hence to make second term  $\leq 2\varepsilon$ , enough if

$$H \ge H_{\gamma,\varepsilon/\sqrt{d}} \quad \text{and} \quad m \ge \frac{d}{(1-\gamma)^2 \varepsilon^2} \log \frac{d(d+1)}{\delta}$$
  
Total # samples:  $|\mathcal{C}| Hm \approx \frac{d^3 \log(d/\varepsilon) \log(d/\delta)}{(1-\gamma)^3 \varepsilon^2}$ 

# Questions from slack



Generally, for two uncorrelated random variables X and Y, the magnitude of the sum of XY over some samples could be much smaller than the sum of |X| |Y|. Now in the proof if  $\varepsilon$  is mostly made up of a "variance" component that is not correlated with the features, can we reduce its effect on the bound in this way? Would the improvement in the bound worth the extra hassle?

And does the "helpful averaging" in the notes refer to this?!

+8

Go ((2))(2)) - + (G

## Discussion



**Theorem (Runtime Bound for Policy Iteration):** Consider a finite, discounted MDP with rewards in [0, 1]. Let  $k^*$  be as in the progress lemma,  $\{\pi_k\}_{k\geq 0}$  the sequence of policies obtained by policy iteration starting from an arbitrary initial policy  $\pi_0$ . Then, after at most  $k = k^*(SA - S) = \tilde{O}\left(\frac{SA - S}{1 - \gamma}\right)$  iterations, the policy  $\pi_k$  produced by policy iteration is optimal:  $v^{\pi_k} = v^*$ . In particular, policy iteration computes an optimal policy with at most  $O\left(\frac{S^4A + S^3A^2}{1 - \gamma}\right)$  arithmetic and logic operations.

Is this satisfactory? Does this mean that if PI is coded up, it will give the optimal policy?

#### Problems

No infinite precision arithmetic on computers Floating point is funky:

Errors can propagate, get large, overwhelm, ..

Order of operations matter sum([1]+[1.0/1 billion]\*1 billion) == 2?

Can't invert "ill conditioned" matrices. Do we have those?

## The goal

# To know whether some calculations are "safe" **Proposition:**

With floating point using e + f bits and target accuracy  $\varepsilon > 0$ , effective horizon  $H = 1/(1 - \gamma)$ , provided  $SA \le u(e, f, 1/\varepsilon)$ , for any MDP M of this size, policy iteration returns an  $\varepsilon$ -optimal policy.

What is u? How do we get u?



# Approaches: Computation over reals

Want computation over the "reals"

- E.g. cost of  $A \mapsto A^{-1}b$ , or  $x \mapsto \sqrt{x}$
- Model 1: Turing machines
- Model 2: #operations in program that uses infinite precision arithmetic (BSS model)
- Model 3: Bit model. Program runs on Turing machine, gets as input target accuracy, computes desired precision of input, gets arbitrarily rounded inputs

What are the strengths and weaknesses?