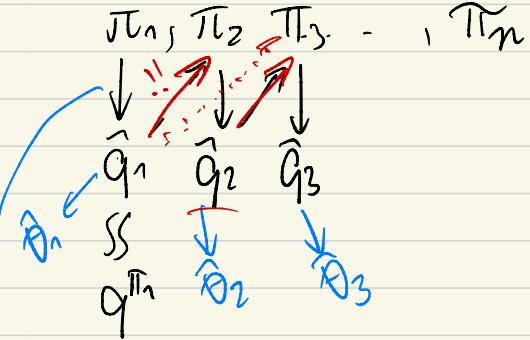


March 4

## POLITEX

## Policy Iteration with Expert Advice



$$\tilde{\theta}_{k-1}^T \psi(s, a)$$

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

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

LSPE G-optimal design  
m rollouts, H length,

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

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

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

Output:  $\pi = \frac{1}{n} (\pi_1 + \dots + \pi_n)$ .

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

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

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

$I \leq f_n$

Online Learning = Optimization

$$(1-\gamma) \sqrt{n}$$

$$\|\mathbb{D}_*\| \leq \frac{2}{1-\gamma} \max_{1 \leq k \leq n} \|q^{\pi_k} - \hat{q}_k\|_\infty \leq \boxed{\frac{2\epsilon(1+\sqrt{d})}{1-\gamma}} + \boxed{\frac{2\sqrt{d}}{(1-\gamma)^2} \left( \gamma^H + \sqrt{\frac{\log((d+1)/\delta)}{2m}} \right)}$$

**Lemma (LSPE-G extrapolation error control):** Fix any full-rank feature-map  $\varphi : \mathcal{Z} \rightarrow \mathbb{R}^d$  and take the set  $\mathcal{C} \subset \mathcal{Z}$  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. (5). Then, for any  $0 \leq \delta \leq 1$ , with probability  $1 - \delta$ ,

$$\|q^\pi - \Phi\hat{\theta}\|_\infty \leq \|\varepsilon_\pi\|_\infty (1 + \sqrt{d}) + \sqrt{d} \left( \frac{\gamma^H}{1-\gamma} + \frac{1}{1-\gamma} \sqrt{\frac{\log(2|\mathcal{C}|/\delta)}{2m}} \right). \quad (6)$$

$E[\mathbb{D}_*] \leq \mathcal{S} = \boxed{\frac{\sigma_0}{n}}$ , truncate  $\hat{q}_k$  to  $[0, \frac{1}{1-\gamma}]$

$\mathcal{S} \geq \frac{1}{1-\gamma} + \dots$  assuming  $\Gamma_\alpha(\varrho) \in [0, 1]$

---

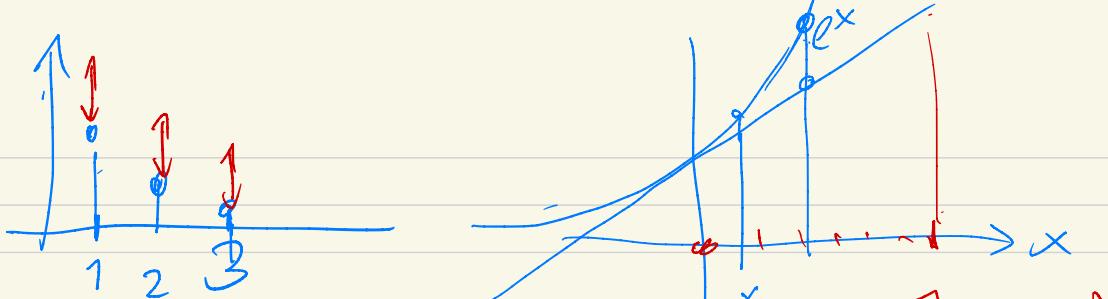
$\mathbb{D}_* = \sum_{k=1}^n M_{\pi^*} \hat{q}_k - M_{\pi_k} \hat{q}_k$  | ses  
 $I(s) = \sum_{k=1}^n \hat{q}_k(s, \pi^*(s)) - \sum_{a \in \mathcal{A}(s)} \pi_k(a | s) \hat{q}_k(s, a)$

Choose  $\pi_k \propto \exp(\gamma \sum_{j=1}^{k-1} q_j)$

1.  $\frac{1}{1-\gamma} \sqrt{2n \log \alpha}$   
2. Compute  $\hat{q}_k$  [depends on  $\pi_k$ ]  
 $\hat{q}_k \in [0, \frac{1}{1-\gamma}]$

$$\pi_k \propto \exp(\gamma \hat{q}_{k-1}) \quad \gamma \rightarrow \infty$$

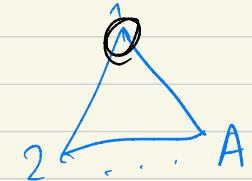
$$\pi_k(a | s) \propto \exp(\gamma (\hat{q}_{k-1}(s, a) + c)), \quad a \in \mathcal{A}$$



$$\frac{\exp(\gamma \hat{q}_{k-1}(s)) e^{\gamma c}}{\sum_a \exp(\gamma \hat{q}_{k-1}(s)) e^{\gamma c}}$$

Fixing  $s \in S$ .

$k = 1, \dots, n$ :



① Learner chooses  $\pi_k \in M_1(\mathcal{V})$

② Adversary  $\hat{q}_k: \mathcal{V} \rightarrow [0, \frac{1}{1-\delta}]$  (reactive)

+  $\boxed{\pi^* \in M_1(\mathcal{A})}$

Need to control:  $R_n = \sum_{k=1}^n \langle \pi^*, \hat{q}_k \rangle - \langle \pi_k, \hat{q}_k \rangle$

$$R_n = O(n)$$

$$\frac{R_n}{n} \rightarrow 0, n \rightarrow \infty$$

repeat

"Online linear optimization"

$\subseteq$  "Online learning"

$$\boxed{\pi \mapsto \langle \pi, \hat{q}_k \rangle}$$

$l_1, \dots, l_n : \mathcal{X} \rightarrow \mathbb{R}$ ,  $\mathcal{X} \subseteq \mathbb{R}^d$   
 $x^* \in \mathcal{X}$   
 $\mathcal{X} \neq \emptyset$   
 convex (closed)

$$R_n = \sum_{t=1}^n l_t(x) - \sum_{t=1}^n l_t(x^*)$$

if  $l \in \text{convex}$

1. Convex opt.

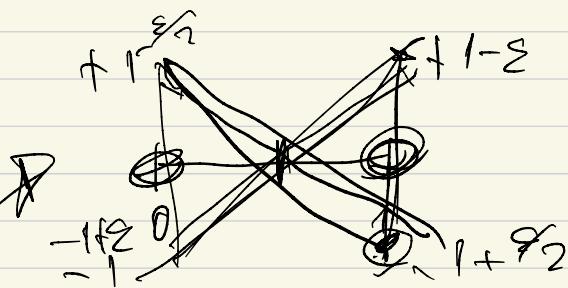
$$l_1 = l_2 = \dots = l_n = l$$

2. Gradient

3. Nonconvex opt.

$$l_t(x) = \langle x_t, y_t \rangle$$

4. Learning



Theorem:  $\pi_k(a) \propto \exp(\eta \sum_{j=1}^{k-1} \hat{q}_j(\cdot))$

$$R_n \leq \frac{1}{1-\gamma} \sqrt{2n \log(A)} \quad \rightarrow \gamma = ?$$

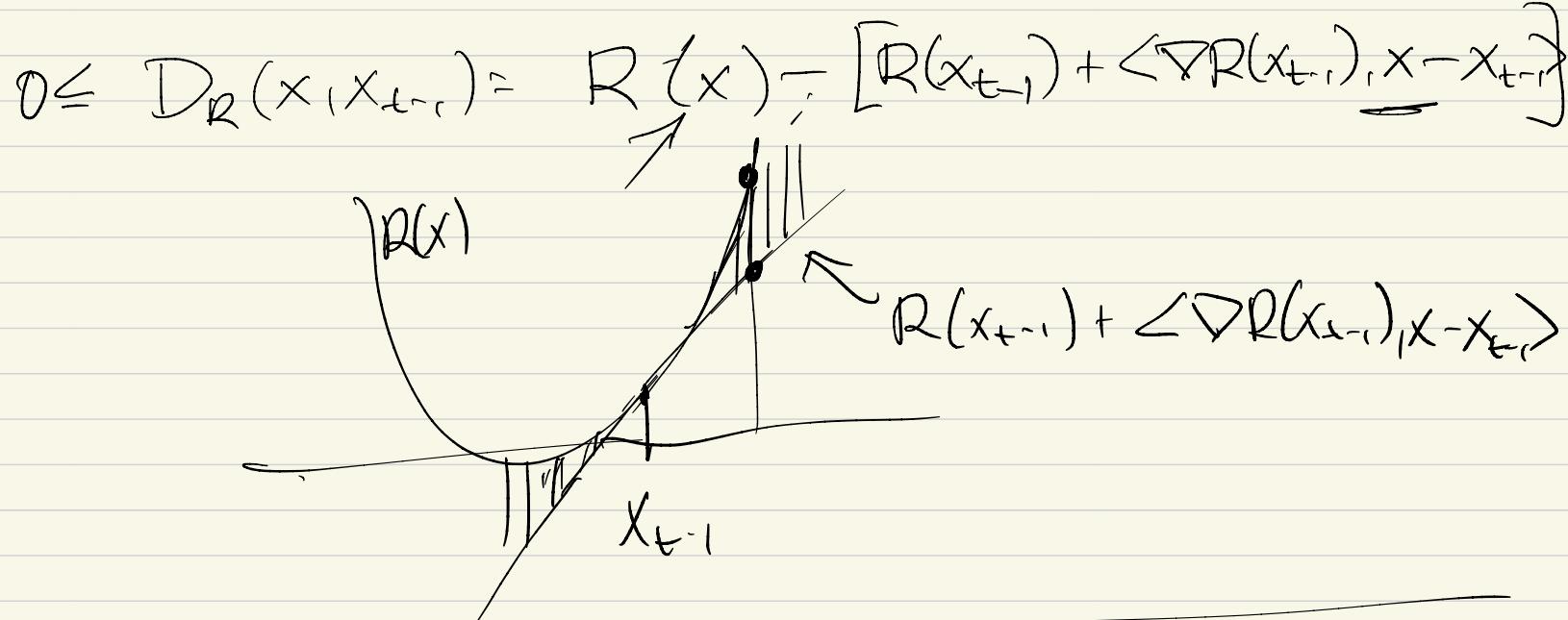
$$\gamma = \sqrt{\frac{2 \log A}{n}}$$

FTL :  $x_t = \underset{x \in \mathcal{X}}{\operatorname{argmin}} \sum_{s=1}^{t-1} l_s(x)$

FTRL :  $x_t = \underset{x \in \mathcal{X}}{\operatorname{argmin}} \sum_{s=1}^{t-1} l_s(x) + \eta R(x)$

convex  $R: \mathcal{X} \rightarrow \mathbb{R}$   $\eta > 0$

$$\text{MD} : x_t = \underset{x \in X}{\operatorname{argmin}} \quad l_{t-1}(x) + \gamma D_R(x, x_{t-1})$$



$$X = M_1(\mathbb{M}) = \{ p \in [0,1]^A \mid \sum_{i=1}^A p_i = 1 \}$$

$$R(p) = - \sum_{i=1}^A p_i \log p_i \quad | \quad R(p) = \frac{1}{2} \|p\|_2^2$$

Tsallis entropy

$$D_R(p, q) = KL(p, q)$$

$\sqrt{A}$

$$p_t = \underset{p \in M_1(A)}{\operatorname{argmin}} \langle p, y_{t-1} \rangle + \gamma KL(p, p_{t-1})$$

D

$$p_t(a) \propto \underbrace{p_{t-1}(a)} e^{-\gamma y_{t-1}(a)}, \quad a \in A$$

+ ~~convex~~

$$P_1(a) = \frac{1}{A}$$

$$-\gamma(y_{t+1}(a) + \dots + y_1(a))$$

$$P_t(a) \propto e$$