

Contextual Bandits w/ policies  $\Pi$ :  $\pi \in \Pi$  is a map  $\bar{\pi}: \mathcal{C} \rightarrow \mathcal{X}$

Input  $\Pi$

for  $t=1, 2, \dots, T$

Nature reveals context  $c_t \in \mathcal{C}$  is iid from  $\mathcal{D}$

Player choose action  $x_t \in \mathcal{X}$  (equiv. chooses  $\pi_t \in \Pi$  and plays  $\pi_t(c_t)$ )

Player receives  $r_t \in [0, 1]$  where  $\mathbb{E}[r_t | c_t, x_t] = v(c_t, x_t)$

Define  $V(\pi) = \mathbb{E}_c [v(c, \pi(c))]$

$$\text{Regret} = \max_{\pi} \sum_{t=1}^T V(\bar{\pi}) - V(\pi_t) \quad (1)$$

$$\max_{\pi \in \Pi} \sum_{t=1}^T v(c_t, \pi) - v(c_t, \pi_t) \quad (2)$$

$$\max_{\pi} \sum_{t=1}^T r_t[\pi(c_t)] - r_t[\pi_t(c_t)] \quad (3)$$

In our stochastic setting these are all within  $O(\sqrt{T})$

Idea: play some logging policy  $\mu$ : in response

to context  $c_t$ , policy plays  $x_t \in \mathcal{X}$  w.p.  $\mu(x | c_t)$

Collect dataset  $\{(c_t, x_t, r_t, p_t)\}_{t=1}^T$   $p_t := \mu(x_t | c_t)$

$$\text{Estimate } \hat{V}(\pi) = \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{1}\{x_t = x\}}{p_t} r_t$$

$$\mathbb{E}[\hat{V}(\pi)] = V(\pi)$$

$$\mathbb{E}[(\hat{V}(\pi) - V(\pi))^2] \leq \mathbb{E}_{c \sim \mathcal{D}} \left[ \frac{1}{\mu(\pi(c) | c)} \right].$$

By Bernstein. w.p.  $\geq 1-\delta$ ,  $\forall \pi \in \Pi$  we have

$$|\hat{V}(\pi) - V(\pi)| \leq \sqrt{\mathbb{E}_{c \sim \mathcal{D}} \left[ \frac{1}{\mu(\pi(c)|c)} \right]} \cdot \frac{2 \log(2/\delta)}{3} + \frac{2 V_{\max} \log(2/\delta)}{3}$$

where  $V_{\max} = \max_{c, x} \frac{1}{\mu(x, c)}$ .

$\mathcal{I}$ -Greedy: Algorithm for "model the bias"/general policies

Input  $\Pi$

for  $t=1, 2, \dots, \mathcal{I}$

Nature reveals  $c_t$

Alg plays  $x_t$  uniformly at random

Construct  $\hat{V}(\pi)$  as above

$$\hat{\pi} = \operatorname{argmax}_{\pi} \hat{V}(\pi)$$

for  $t = \mathcal{I}+1, \dots, T$

Nature reveals  $c_t$

$$x_t = \hat{\pi}(c_t)$$

$$\mu(x|c_t) = \frac{1}{|\mathcal{X}|} \text{ for all } x \in \mathcal{X}.$$

$$\Rightarrow \text{w.p. } \geq 1-\delta \quad \forall \pi \in \Pi \quad |\hat{V}(\pi) - V(\pi)| \leq \underbrace{\sqrt{\frac{4|\mathcal{X}| \log(2/\delta)}{3}}}_{:= \epsilon_{\mathcal{I}}}$$

$$\pi_{\star} \in \operatorname{argmax}_{\pi \in \Pi} V(\pi)$$

$$\begin{aligned}
 V(\hat{\pi}) &= \underbrace{V(\hat{\pi}) - \hat{V}(\hat{\pi})}_{\geq -\varepsilon_{\mathcal{I}}} + \underbrace{\hat{V}(\hat{\pi}) - \hat{V}(\pi_*)}_{\geq 0} + \underbrace{\hat{V}(\pi_*) - V(\pi_*)}_{\geq -\varepsilon_{\mathcal{I}}} + V(\pi_*) \\
 &\geq V(\pi_*) - 2\varepsilon_{\mathcal{I}}
 \end{aligned}$$

$$\begin{aligned}
 \text{Regret} &= \sum_{t=1}^T V(\pi_*) - V(\pi_t) \\
 &= \sum_{t=1}^{\mathcal{I}} V(\pi_*) - V(\pi_t) + \sum_{t=\mathcal{I}+1}^T V(\pi_*) - V(\pi_t) \\
 &\leq \mathcal{I} \cdot 1 + (T-\mathcal{I}) \underbrace{(V(\pi_*) - V(\hat{\pi}))}_{\leq 2\varepsilon_{\mathcal{I}}} \\
 &= \mathcal{I} + 2(T-\mathcal{I}) \sqrt{\frac{4|\chi| \log(2/\pi/\delta)}{\mathcal{I}}} \\
 &= O\left(T^{2/3} (|\chi| \log(2/\pi/\delta))^{1/3}\right)
 \end{aligned}$$

$$\mathcal{I} = (|\chi| T^2 \log(2/\pi/\delta))^{1/3}$$

How do we find  $\operatorname{argmax}_{\pi \in \Pi} \hat{V}(\pi)$  efficiently?

$$\begin{aligned} \hat{V}(\pi) &= \frac{1}{3} \sum_{t=1}^3 \frac{\mathbb{1}\{x_t = \pi(c_t)\}}{P_t} r_t \\ &= \frac{1}{3} \sum_{t=1}^3 \frac{(1 - \mathbb{1}\{x_t \neq \pi(c_t)\})}{P_t} r_t \\ &= \frac{1}{3} \sum_{t=1}^3 \frac{r_t}{P_t} - \frac{1}{3} \sum_{t=1}^3 \mathbb{1}\{x_t \neq \pi(c_t)\} \frac{r_t}{P_t} \end{aligned}$$

$$\operatorname{argmax}_{\pi} \hat{V}(\pi) = \operatorname{argmin}_{\pi} \frac{1}{3} \sum_{t=1}^3 \mathbb{1}\{x_t \neq \pi(c_t)\} \frac{r_t}{P_t}$$

Think of each  $\pi: \mathcal{C} \rightarrow \mathcal{X}$  as a classifier over classes  $\{1, \dots, |\mathcal{X}|\}$  and examples in  $\mathcal{C}$ .

$x_t \in \mathcal{X}$  is a "label",  $c_t \in \mathcal{C}$  example,  $\frac{r_t}{P_t}$  is weight

Minimize 0/1 loss on  $\{(c_t, x_t, \frac{r_t}{P_t})\}_{t=1}^3$   
 $\uparrow$   
 weight of  $t$ th point

Ex. let  $f_{\theta}: \mathcal{C} \rightarrow \mathbb{R}^{\mathcal{X}}$  be a neural network parameterized by parameters  $\theta \in \mathbb{R}^d$

$$\text{loss}(\theta) = -\frac{1}{3} \sum_{t=1}^3 \frac{r_t}{P_t} \log \left( \frac{\exp(f_{\theta}(c_t)[x_t])}{\sum_{x \in \mathcal{X}} \exp(f_{\theta}(c_t)[x])} \right)$$

$$\text{Let } \hat{\theta}_3 = \underset{\theta}{\text{argmin}} \text{Loss}(\theta)$$

When  $c_t$  arrives we play

$$x_t = \underset{x \in \mathcal{X}}{\text{argmax}} \left[ f_{\hat{\theta}_3}(c_t) \right]_x$$

Zooming out:  $\gamma$ -greedy achieves  $R_T \leq T^{2/3} (|\mathcal{X}| \log(|\mathcal{T}|/\delta))^{1/3}$ .

But  $|\mathcal{C}|=1$  and  $|\mathcal{T}|=|\mathcal{X}|$  in the standard MAB then we know  $\sqrt{|\mathcal{X}|T}$  regret is possible.

→ For contextual Bandits, can we achieve  $R_T \leq \sqrt{|\mathcal{X}|T \cdot \log(|\mathcal{T}|)}$ ?

Answer: yes.

An elimination alg for stochastic Contextual Bandits

Input  $\mathcal{T}, \delta$

$$\hat{\mathcal{T}}_1 = \mathcal{T}$$

for  $l=1, 2, \dots$

$$\mathcal{E}_l = \zeta^{-l}, \quad \begin{array}{ll} \# \text{ measurements} & \mathcal{J}_l \text{ TBD} \\ \text{regularization} & \gamma_l \text{ TBD} \end{array}$$

$$Q_l = \underset{Q \in \Delta_{\hat{\mathcal{T}}_l}}{\text{argmin}} \max_{\pi \in \hat{\mathcal{T}}_l} \mathbb{E}_c \left[ \frac{1}{Q^x(\pi(c)|c)} \right]$$

$$\text{where } Q^x(x|c) := \gamma + (1-\gamma)\delta Q(x|c)$$

$$Q(x|c) := \sum_{\pi: \pi(c)=x} Q(\pi) \leftarrow \text{prob dist over actions}$$

$$T_\epsilon = \sum_{i=1}^{\frac{1}{\epsilon}} \mathcal{I}_\epsilon$$

for  $t = T_{\epsilon-1} + 1, \dots, T_\epsilon$

Nature reveals  $C_t$

Player plays  $x_t \sim Q_\epsilon^\gamma(\cdot | C_t)$  to get reward  $r_t$

$$\hat{V}_\epsilon(\pi) = \frac{1}{\mathcal{I}_\epsilon} \sum_{t=1}^{\mathcal{I}_\epsilon} \frac{\mathbb{1}\{\pi(C_t) = x_t\}}{P_t} r_t$$

$$\hat{\Pi}_{\epsilon+1} = \hat{\Pi}_\epsilon \setminus \{ \pi \in \hat{\Pi}_\epsilon : \max_{\pi'} \hat{V}(\pi') - \hat{V}(\pi) \geq 2\epsilon \}$$

//  $Q_\epsilon$  is chosen to minimize maximum variance of  $\hat{V}(\pi)$ :

$$Q_\epsilon \approx \underset{Q \in \Delta_{\mathcal{X}_\epsilon}}{\operatorname{argmin}} \max_{\pi \in \hat{\Pi}_\epsilon} \mathbb{E} [ (\hat{V}(\pi) - V(\pi))^2 ]$$

Lemma For any policy  $\Pi$  and dist. over contexts

$$\min_{Q \in \Delta_{\mathcal{X}_\epsilon}} \max_{\pi \in \Pi} \mathbb{E} \left[ \frac{1}{Q(\pi(c) | c)} \right] \leq |\mathcal{X}|.$$

Furthermore, if  $\gamma < \frac{1}{2K}$  then

$$\min_{Q \in \Delta_{\mathcal{X}_\epsilon}} \max_{\pi \in \Pi} \mathbb{E} \left[ \frac{1}{Q^\gamma(\pi(c) | c)} \right] \leq 2|\mathcal{X}|.$$

This  $\nearrow$  lemma is a corollary of Kiefer-Wolfowitz which we will show next time.

Consequently we have by Bernstein

$$|\hat{V}_\varepsilon(\bar{\tau}) - V(\bar{\tau})| \leq \sqrt{\frac{(2|\chi|) \cdot 2 \log(2|\pi|/\delta)}{3\varepsilon}} + \frac{2 \log(2|\pi|/\delta)}{\delta \cdot 3\varepsilon}$$

$$\leq \sqrt{\frac{16|\chi| \log(2|\pi|/\delta)}{3\varepsilon}} =: \varepsilon_\varepsilon$$

$$\text{for } \delta_\varepsilon = \min \left\{ \frac{1}{2|\chi|}, \sqrt{\frac{2 \log(2|\pi|/\delta)}{9|\chi| \cdot 3\varepsilon}} \right\}$$