

Contextual Bandits - "Model the world"

Input  $\mathcal{F}$ ,  $f \in \mathcal{F}$   $f: \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}$ ,  $|\mathcal{A}| = n$

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

Adversary chooses  $c_t \in \mathcal{C}$  arbitrarily

Player chooses  $a_t \in [n]$

Nature reveals  $y_t \in [0, 1]$ ,  $\mathbb{E}[y_t | c_t, a_t] = f_{\#}(c_t, a_t)$

minimize Regret =  $\sum_{t=1}^T \max_a f_{\#}(c_t, a) - f_{\#}(c_t, a_t)$

Assume realizability:  $f_{\#} \in \mathcal{F}$ .

UCB: Assume  $\exists \phi: \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}^d$ ,  $\mathcal{F} = \{ \langle \theta, \phi(c, a) \rangle, \theta \in \mathbb{R}^d \}$

$$a_t = \underset{a}{\operatorname{argmax}} \max_{\theta \in \mathcal{C}_t} \langle \theta, \phi(c_t, a) \rangle$$

$$\mathcal{C}_t = \{ \theta : \|\theta - \hat{\theta}_t\|_{V_t}^2 \leq \beta_t^2 \}$$

We showed that for any sequence of  $x_1, x_2, \dots$

and  $y_t = \langle \theta, x_t \rangle + \eta_t$  we have  $\|\theta - \hat{\theta}_t\|_{V_t}^2 \leq \beta_t^2$

w.p.  $\geq 1 - \delta$  where  $V_t = \lambda I + \sum_s x_s x_s^T$ ,  $\hat{\theta}_t = V_t^{-1} \sum_s x_s y_s$

to apply same result w/  $x_t = \phi(c_t, a_t)$ .

$$\Rightarrow \text{Regret} \leq d\sqrt{T}.$$

Thompson Sampling

$$V_t = \lambda I + \sum_{s=1}^t \phi(c_s, a_s) \phi(c_s, a_s)^T$$

$$\hat{\theta}_t = V_t^{-1} \sum_{s=1}^t \phi(c_s, a_s) y_s$$

At time  $t$ , draw

$$\tilde{\theta}_t \sim \mathcal{N}(\hat{\theta}_t, \alpha V_t^{-1})$$

$$\text{Play } a_t = \underset{a \in \{1, \dots, n\}}{\text{argmax}} \phi(c_t, a)^T \tilde{\theta}_t$$

Greedy

$$\text{Play } a_t = \underset{a}{\text{argmax}} \langle \hat{\theta}_t, \phi(c_t, a) \rangle$$

Can work! If  $\min_a \sum_t \phi(c_t, a) \phi(c_t, a)^T \succ 0$

Thompson Sampling is great. How do I generalize it to arbitrary  $\mathcal{F}$ ?

Note for linear

$$\bar{\theta}_t = \left( \sum_{s=1}^t x_s x_s^T \right)^{-1} \sum_{s=1}^t x_s (y_s + z_s'), \quad z_s' \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$$

$$= \hat{\theta}_t + \underbrace{V_t^{-1} \sum_{s=1}^t x_s z_s'}_{\sim \mathcal{N}(0, V_t^{-1} (\sum_{s=1}^t x_s x_s^T) V_t^{-1})}$$

$$\sim \mathcal{N}(0, V_t^{-1} (\underbrace{\sum_{s=1}^t x_s x_s^T}_{V_t}) V_t^{-1})$$

$$= \mathcal{N}(0, V_t^{-1})$$

$$\bar{\theta}_t \sim \mathcal{N}(\hat{\theta}_t, V_t^{-1})$$

Generalize Thompson Sampling:

$$\bar{f}_t = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{s=1}^t (y_s + z_s' - f(l_s, a_s))^2$$

$$\text{Play } \operatorname{argmax}_a \bar{f}_t(l_t, a).$$

How do we generalize optimism to arbitrary function classes?

Suppose I collected data  $\{(c_s, a_s, y_s)\}_t$

~~$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{H}} \sum_t (f(c_t, a_t) - y_t)^2$$~~

$$\hat{f}_t \leftarrow \text{ORO}(\{(c_s, a_s, y_s)\}_{s=1}^{t-1})$$

$$a_t = \operatorname{argmax}_a \max_{f \in \mathcal{H}_t} f_+(c_t, a)$$

where

$$\mathcal{H}_t = \left\{ f : \sum_{s=1}^t \mathbb{E}_{a_s \sim p_s} \left[ (\hat{f}_s(c_s, a_s) - f(c_s, a_s))^2 \right] \leq \log |\mathcal{X} / \delta| \right\}$$

Generalized UCB.  $\uparrow$

Online Regression Oracle At each time  $t$

we have  $\{(c_s, a_s, y_s)\}_{s=1}^{t-1}$  w/  $\mathbb{E}[y_s | c_s, a_s] = f_+(c_s, a_s)$

and  $a_s \sim p_s \in \Delta_n$ . An ORO returns  $\hat{f}_t$ :

$$\sum_{t=1}^T \mathbb{E}_{a_t \sim p_t} \left[ (\hat{f}_t(c_t, a_t) - f_a(c_t, a_t))^2 \right] \leq \text{Est}(\mathcal{F}, \delta, T)$$

w.p.  $\geq 1 - \delta$ .

Theorem] Play Exponential weights over  $\mathcal{F}$  with

loss  $l_{t,i} = (f^i(c_t, a_t) - y_t)^2$  and  $\xi = 1/2$

and output  $\hat{f}_t(c, a) = \sum_i f^i(c, a) p_{t,i}$  then

$$\text{Est}(\mathcal{F}, \delta, T) \leq \log(|\mathcal{F}|/\delta)$$

Theorem]  $\exists$   $\mathcal{F}$  and distribution of contexts such that

Generalized UCB suffers linear regret

Proof  $\mathcal{C} = [m]$   $\mathcal{A} = \{1, 2\}$  Fix  $\epsilon > 0$

$$f_a(c, 1) = 1 - \epsilon$$

$$f_a(c, 2) = 0$$

$c_t \sim \text{Uniform}([m])$

$$y_t = f_a(c_t, a_t)$$

$$f_i(c_i, 1) = 0$$

$$f_i(c_i, 2) = 1 \quad \forall c_i$$

$$\Rightarrow \text{Regret} \geq \min\{|\mathcal{C}|, |\mathcal{F}|\}$$

Optimism can be a very bad idea.

$\epsilon$ -Greedy

Input:  $\epsilon > 0$

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

$$P_t(a) = \frac{\epsilon}{A} + \mathbb{1}\{\hat{a}_t = a\} (1-\epsilon)$$

Obtain  $\vec{f}_t$  from  $\text{ORO}(\{(c_s, a_s, y_s)_{s \leq t}\})$

Observe context  $c_t$

w/ prob  $\epsilon$  play  $a_t \sim \text{uniform}(\{n\})$

w/ prob  $1-\epsilon$  play  $a_t = \hat{a}_t = \underset{a}{\text{argmax}} \vec{f}_t(c_t, a)$

Observe reward  $y_t \in [0, 1]$

Theorem | Assume  $\text{ORO}$  satisfies  $\text{Est}(\mathcal{F}, T, \delta) \leq \log(|\mathcal{F}|/\delta)$ .

Then  $\epsilon$ -greedy satisfies

$$\sum_{t=1}^T f_{\star}(c_t, a_{\star}(c_t)) - f_{\star}(c_t, a_t) \leq n^{1/3} T^{2/3} \log(|\mathcal{F}|/\delta)^{1/3}$$

w.p.  $\geq 1-\delta$ .

Proof

$$\sum_{t=1}^T f_{\star}(c_t, a_{\star}(c_t)) - f_{\star}(c_t, a_t)$$

$$\leq \epsilon T + \sum_{t=1}^T f_{\star}(c_t, a_{\star}(c_t)) - f_{\star}(c_t, \hat{a}_t)$$

$$a_t := a_t(t)$$

$$f_t(t, a_t) - f_t(t, \hat{a}_t) = f_t(t, a_t) - \hat{f}(t, a_t)$$

$$+ \hat{f}(t, a_t) - \hat{f}(t, \hat{a}_t) \stackrel{\leq 0}{\leq 0}$$

$$P_t(a) = \frac{\varepsilon}{A} + \mathbb{1}\{\hat{a}_t = a\} (1 - \varepsilon)$$

$$+ \hat{f}(t, \hat{a}_t) - f_t(t, \hat{a}_t)$$

$$\leq \sum_{a \in \hat{a}_t, a_t} |f_t(t, a) - \hat{f}(t, a)|$$

$$= \sum_{a \in \hat{a}_t, a_t} \frac{1}{\sqrt{P_t(a)}} \cdot \sqrt{P_t(a)} |f_t(t, a) - \hat{f}(t, a)|$$

$$\leq \left( \frac{2n}{\varepsilon} \right)^{1/2} \left( \sum_{a \in \hat{a}_t, a_t} P_t(a) (f_t(t, a) - \hat{f}(t, a))^2 \right)^{1/2}$$

$$\leq \left( \frac{2n}{\varepsilon} \right)^{1/2} \mathbb{E}_{a \sim P_t} \left[ (f_t(t, a) - \hat{f}(t, a))^2 \right]^{1/2}$$

$$\text{Regret}_T \leq \varepsilon T + \sum_{t=1}^T \sqrt{\frac{2n}{\varepsilon} \mathbb{E}_{a \sim P_t} \left[ (f_t(t, a) - \hat{f}(t, a))^2 \right]}$$

$$\leq \varepsilon T + \sqrt{\frac{2nT}{\varepsilon} \log(|\mathcal{F}|/\delta)}$$

$$\leq T^{2/3} \cdot (2n \log(|\mathcal{F}|/\delta))^{1/3} \quad \text{for some } \varepsilon.$$

## Inverse Gap Weighting

Fix  $\gamma > 0$  and  $g: [n] \rightarrow \mathbb{R}$ , let  $\hat{a} = \operatorname{argmax}_a g(a)$ .

Define  $\text{IGW}_\gamma(a) := p(a) = \frac{1}{\lambda + 2\gamma(g(\hat{a}) - g(a))}$

where  $\lambda$  is chosen s.t.  $\sum_a p(a) = 1$ .

$$\frac{1}{\lambda} \leq \sum_a p(a) \leq \frac{n}{\lambda} \Rightarrow \lambda \in [1, n]$$

Lemma Fix some  $g: [n] \rightarrow \mathbb{R}$  and  $p = \text{IGW}_\gamma$

Then for any  $g_a: [n] \rightarrow \mathbb{R}$  we have

$$\mathbb{E}_{a \sim p} [g_a(a_a) - g_a(a)] \leq \frac{n}{\gamma} + \gamma \mathbb{E}_{a \sim p} [(g(a) - g_a(a))^2]$$

where  $a_a = \operatorname{argmax}_a g_a(a)$ .

Proof

(I)

(II)

$$\mathbb{E}_{a \sim p} [g_a(a_a) - g_a(a)] = \mathbb{E}_{a \sim p} [g(\hat{a}) - g(a)] + \mathbb{E}_{a \sim p} [g(a) - g_a(a)]$$

$$g_a(a_a) - g(\hat{a}) \quad \text{(III)}$$

$$\begin{aligned}
 \text{(I)} \quad \mathbb{E}_{a \sim p} [g(\hat{a}) - g(a)] &= \sum_a p(a) \cdot (g(\hat{a}) - g(a)) \\
 &= \sum_i \frac{(g(\hat{a}_i) - g(a_i))}{\lambda + \gamma (g(\hat{a}_i) - g(a_i))} \\
 &\leq \frac{\eta}{2\gamma}
 \end{aligned}$$

$$\text{(II)} \quad u \leq \frac{1}{2\gamma} + \frac{\gamma}{2} u^2$$

$$\mathbb{E}_{a \sim p} [g(a) - g_{\#}(a)] \leq \frac{1}{2\gamma} + \frac{\gamma}{2} \mathbb{E}_{a \sim p} [(g(a) - g_{\#}(a))^2]$$

$$\text{(III)} \quad u \leq \frac{1}{2\gamma p(a_{\#})} + \frac{\gamma p(a_{\#})}{2} u^2$$

$$\begin{aligned}
 g_{\#}(a_{\#}) - g(a_{\#}) \\
 \leq \frac{1}{2\gamma p(a_{\#})} + \frac{\gamma}{2} p(a_{\#}) (g_{\#}(a_{\#}) - g(a_{\#}))^2
 \end{aligned}$$

$$p(a_n) = \frac{1}{\lambda + 2\gamma (g(\hat{a}) - g(a_n))} \geq$$

$$\frac{1}{2\gamma p(a_n)} = \frac{\lambda + 2\gamma (g(\hat{a}) - g(a_n))}{2\gamma}$$

$$g(a_n) - g(\hat{a})$$

$$= g_a(a_n) - g(a_n) + g(a_n) - g(\hat{a})$$

$$\leq \frac{1}{2\gamma p(a_n)} + \frac{\gamma}{2} p(a_n) (g_a(a_n) - g(a_n))^2 + g(a_n) - g(\hat{a})$$

$$= \frac{\lambda + 2\gamma (g(\hat{a}) - g(a_n))}{2\gamma} + \frac{\gamma}{2} p(a_n) (g_a(a_n) - g(a_n))^2 + g(a_n) - g(\hat{a})$$

$$\leq \frac{n}{2\gamma} + \frac{\gamma}{2} \mathbb{E}_{a \sim p} [(g_a(a) - g(a))^2]$$

# Square CB

Input  $\gamma > 0$

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

get  $\hat{f}_t$  from  $\text{ORO}(\{C_s, a_s, y_s\}_{s \leq t})$

Observe  $C_t$ , let  $\hat{a}_t = \underset{a}{\operatorname{argmax}} \hat{f}_t(C_t, a)$

Play  $a_t \sim P_t(C_t, a) = \frac{1}{\lambda + 2\gamma(\hat{f}_t(C_t, \hat{a}_t) - \hat{f}_t(C_t, a))}$

Observe  $y_t$

Theorem) Square CB enjoys regret  $\frac{nT}{\gamma} + \gamma \log(|\mathcal{F}|/\delta)$

w.p.  $\geq 1 - \delta$ .

$$\sum_t \mathbb{E}_{a \sim P_t} [f_*(C_t, a_t(C_t)) - f_*(C_t, a)]$$

$$\leq \frac{nT}{\gamma} + \gamma \sum_{t=1}^T \mathbb{E}_{a \sim P_t} [(f_*(C_t, a) - \hat{f}_t(C_t, a))^2]$$

$$\leq \frac{nT}{\gamma} + \gamma \log(|\mathcal{F}|/\delta) \leq \sqrt{nT \log(|\mathcal{F}|/\delta)}.$$