

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

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

Assume realizability: $f_{\theta_t} \in \mathcal{F}$

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

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

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 :

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

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

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

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

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

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

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

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

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

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

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

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