

Contextual Bandits

Input: Π , $\pi \in \Pi$: $\pi: \mathcal{C} \rightarrow \mathcal{X}$ (context to action)

for $t=1, 2, \dots$

Nature reveals context $c_t \stackrel{iid}{\sim} \mathcal{D}$

Player plays $x_t \in \mathcal{X}$ (equiv. $\pi_t \in \Pi$, $x_t = \pi_t(c_t)$)

and receives reward $r_t \in [0, 1]$ w/ $\mathbb{E}[r_t | c_t, x_t] = v(c_t, x_t)$

Value of policy π defined $V(\pi) := \mathbb{E}_{c \sim \mathcal{D}} [v(c, \pi(c))]$

Off-policy evaluation

Assume dataset collected by a logging policy where

$$x_t \sim \mu(\cdot | c_t) \quad \text{where} \quad p_t = \mu(x_t | c_t)$$

for T steps to construct $\{(c_t, x_t, r_t, p_t)\}_{t=1}^T$

IPS-estimator ("Model the bias")

$$\hat{v}(c, x) = \frac{\mathbb{1}\{x_t = x\}}{p_t} r_t \quad \text{and} \quad \hat{V}(\pi) = \frac{1}{T} \sum_{t=1}^T \hat{v}(c_t, \pi(c_t))$$

We showed $\mathbb{E}[\hat{v}(c, x) | c] = v(c, x) \Rightarrow \mathbb{E}[\hat{V}(\pi)] = V(\pi)$

$$\text{Var}(\hat{v}(c, x) | c) \leq \frac{1}{\mu(x | c)} \quad \text{Var}(\hat{V}(\pi)) \leq \mathbb{E}_{c \sim \mathcal{D}} \left[\frac{1}{\mu(\pi(c) | c)} \right] \cdot \frac{1}{T}$$

Bernstein's Inequality let X_1, \dots, X_n be independent R.V.

w/ $X_i \leq B$, $\text{Var}(X_i) \leq \sigma^2$ then

$$\mathbb{P}\left(\frac{1}{n} \sum X_i - \mathbb{E}[X_i] > \sqrt{\frac{2\sigma^2 \log(1/\delta)}{n}} + \frac{2B \log(1/\delta)}{3n}\right) \leq \delta$$

Recall Hoeffding says if $X_i \in [0, B]$ then

$$\mathbb{P}\left(\frac{1}{n} \sum X_i - \mathbb{E}[X_i] > \sqrt{\frac{B^2 \log(1/\delta)}{2n}}\right) \leq \delta.$$

We always have $r_t \in [0, 1]$

$$\hat{V}(c, x) = \frac{\mathbb{1}\{x_t = x\}}{P_t} r_t \leq \frac{1}{P_t} \leq \max_{c, x} \frac{1}{\mu(x|c)}$$

\Rightarrow For any fixed $\pi \in \Pi$, w.p. $\geq 1 - \delta$.

$$\hat{V}(\pi) - V(\pi) > \sqrt{\mathbb{E}_{c \sim \theta} \left[\frac{1}{\mu(\pi(c)|c)} \right] 2 \log(1/\delta)} \frac{1}{3} + \frac{2 \log(1/\delta)}{33} \max_{c, x} \frac{1}{\mu(x|c)}$$

\Rightarrow w.p. $\geq 1 - \delta$, for all $\pi \in \Pi$ simultaneously

$$\begin{aligned} \text{we have } |\hat{V}(\pi) - V(\pi)| &\leq \sqrt{\mathbb{E}_{c \sim \theta} \left[\frac{1}{\mu(\pi(c)|c)} \right] 2 \log(2|\Pi|/\delta)} \frac{1}{3} \\ &\quad + \max_{c, x} \frac{1}{\mu(x|c)} \cdot \frac{2 \log(2|\Pi|/\delta)}{33} \end{aligned}$$

If $\mu(x|c) = \frac{1}{|\mathcal{X}|}$

$$|\hat{V}(\pi) - V(\pi)| \leq \sqrt{\frac{2|\mathcal{X}| \log(2|\Pi|/\delta)}{3}} + \frac{2|\mathcal{X}| \log(2|\Pi|/\delta)}{33}$$

$$\leq \sqrt{\frac{4|\mathcal{X}| \log(2|\Pi|/\delta)}{3}}$$

\Rightarrow If $3 \geq 4|\mathcal{X}|\bar{\epsilon}^2 \log(2|\Pi|/\delta)$ then
under uniform exploration we have

$$\max_{\pi \in \Pi} |\hat{V}_3(\pi) - V(\pi)| \leq \epsilon \quad \text{w.p.} \geq 1 - \delta.$$

"Model the World"

Fix some function class \mathcal{F} s.t. $f \in \mathcal{F}$

$$f: \mathcal{C} \times \mathcal{X} \rightarrow \mathbb{R}$$

Idea: Ideally $\exists f_* \in \mathcal{F} : v(c, x) \approx f_*(c, x)$

Question: Can we learn f_* w/ our dataset

$$\{(c_t, x_t, r_t, p_t)\}_{t=1}^T$$

$$\hat{f} = \underset{f \in \mathcal{F}}{\text{argmin}} \quad \frac{1}{3} \sum_{t=1}^T (r_t - f(c_t, x_t))^2$$

fix any $f \in \mathcal{F}$

$$\begin{aligned} \mathbb{E}[(r_t - f(c_t, x_t))^2] &= \mathbb{E}[(r_t - v(c_t, x_t) + v(c_t, x_t) - f(c_t, x_t))^2] \\ &= \mathbb{E}[(r_t - v(c_t, x_t))^2] + \mathbb{E}[(v(c_t, x_t) - f(c_t, x_t))^2] \\ &\leq \frac{1}{4} + \sum_{x \in \mathcal{X}} \mathbb{E}[\mathbb{E}[\mathbb{1}\{x_t = x\} (v(c_t, x) - f(c_t, x))^2 | \mathcal{C}_t]] \end{aligned}$$

$$= \frac{1}{4} + \mathbb{E}_c \left[\sum_{x \in \mathcal{X}} \mu(x|c) (v(c,x) - f(c,x))^2 \right]$$

$$\mathbb{E} \left[\frac{(r_t - f(c_t, x_t))^2}{p_t} \right] \leq \frac{1}{4} \mathbb{E} \left[\sum_x \frac{1}{\mu(x|c)} \right] + \mathbb{E} \left[\sum_x (v(c,x) - f(c,x))^2 \right]$$

Note: If $\mu(x|c) > 0 \quad \forall x, c$ and $v \in \mathcal{F}$

then both above objectives satisfy $\hat{f} \rightarrow v$ if $T \rightarrow \infty$.

To evaluate a policy π , output

$$\hat{V}(\pi) = \frac{1}{T} \sum_{t=1}^T \hat{f}(c_t, \pi(c_t))$$

Concerns/Warning:

- Could be wasteful b/c I care about estimating $V(\pi)$, but learning \hat{f} does not take π into account
| \mathcal{F} | too big
- i.e. learning $\hat{f} \in \mathcal{F}$ well enough to estimate $|V(\pi) - \hat{V}(\pi)| \leq \epsilon$ could require way more samples than $\frac{1}{\epsilon^2} \log(|\mathcal{T}|)$
- Biased. If $\min_{f \in \mathcal{F}} \max_{x, c} |f(x, c) - v(c, x)|$ is not small, then regardless of how much
| \mathcal{F} | too small

data you collect, estimates are biased.

Doubly Robust Estimators

Learn some $\hat{f} \in \mathcal{F}$ then set

$$\hat{V}_{DR}(C_t, x) = \hat{f}(C_t, x) + (r_t - \hat{f}(C_t, x)) \frac{\mathbb{1}\{x_t = x\}}{P_t}$$

$$\mathbb{E}[\hat{V}_{DR}(C_t, x) | C_t] = V(C_t, x) \quad \text{Unbiased!}$$

and by computing variance, if $\mathbb{E}[(r_t - V(C_t, x_t))^2]$ is small then variance of $\hat{V}_{DR}(C_t, x)$

$$\mathbb{E}[(V(C_t, x) - \hat{f}(C_t, x))^2] \frac{1}{P_t}$$

This assume $\hat{f} \perp \{(C_t, x_t, r_t, P_t)\}_{t=1}^T$

but in practice \hat{f} is trained on \uparrow

Linear Contextual Bandits (special case: "Model the world")

Assume \exists known $\phi: \mathcal{C} \times \mathcal{X} \rightarrow \mathbb{R}^d$ and unknown $\exists \theta^* \in \mathbb{R}^d$
such that $v(c, x) = \langle \theta^*, \phi(c, x) \rangle$

"Model the world" where $f(c, x) = \langle \theta, \phi(c, x) \rangle$ for $\theta \in \mathbb{R}^d$.

for $t=1, 2, \dots$

Nature reveals $c_t \sim \mathcal{I}$

$$\mathcal{Z}_t = \{ \phi(c_t, x) : x \in \mathcal{X} \}$$

Player chooses $z_t \in \mathcal{Z}_t$

and observes $\langle z_t, \theta^* \rangle + \xi_t$, $\mathbb{E}[\xi_t] = 0$

In practice:

- If \mathcal{X} is unstructured
and $c \in \mathbb{R}^p$, $i \in \mathcal{X}$
 $\phi(c, i) = \text{vec}(c e_i^T)$

- ϕ is learned by
historical data

Natural algorithm is UCB. Construct

confidence set \mathcal{C}_t : $\theta^* \in \mathcal{C}_t \quad \forall t$ w.h.p.

$$\text{Play } z_t = \underset{z \in \mathcal{Z}_t}{\text{argmax}} \max_{\theta \in \mathcal{C}_t} \langle z, \theta \rangle$$

$$R_T \lesssim d \sqrt{T} \quad (\text{ignoring logs})$$

(Thompson sampling works very well here)