

Contextual Bandits (Stochastic)

Input: Π , $\pi_0 \in \Pi$, $\bar{r}: \mathcal{C} \rightarrow [n]$

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

Nature reveals $c_t \stackrel{iid}{\sim} \bar{r}$

Player chooses $\pi_t \in \Pi$, $a_t := \pi_t(c_t)$

Receives reward $r_t \in [0, 1]$: $\mathbb{E}[r_t | c_t] = r(c_t, a_t)$

Minimize regret $\max_{\pi \in \Pi} \sum_{t=1}^T r(c_t, \pi(c_t)) - r(c_t, a_t)$

$$V(\pi) = \mathbb{E}_{c, a \sim \pi(c)} [r(c, a)] \\ = \mathbb{E}_c [r(c, \pi(c))]$$

Logging policy $\mu(\cdot | c_t) \in \Delta_n$, $p_t = \mu(a_t | c_t)$. Assume $\mu(a_t | c_t) > 0 \forall a, c$

Collect dataset $\{(c_t, a_t, r_t, p_t)\}_t$. Always log sampling probs p_t !

Forget context. Given $\{(a_t, r_t, p_t)\}_t$. Idea: learn some function

$$f: f(a_t) \approx r_t \quad (\text{e.g. } \hat{f} = \underset{f \in \mathcal{F}}{\text{argmin}} \sum_t (f(a_t) - r_t)^2)$$

recommend $\underset{a}{\text{argmax}} \hat{f}(a)$

Model the World

Given $\{(c_t, a_t, r_t, p_t)\}_t$

learn $f: f(c_t, a_t) \approx r_t$ (e.g. $\hat{f} = \underset{f \in \mathcal{F}}{\text{argmin}} \sum_t (f(c_t, a_t) - r_t)^2$)

recommend $\underset{a}{\text{argmax}} \hat{f}(c_t, a)$

Model the Bias

$$\hat{r}(c_t, a) := r_t \frac{\mathbb{1}\{a_t = a\}}{P_t}$$

$$\begin{aligned} \mathbb{E}[\hat{r}(c_t, a)] &= \sum_{a'} \frac{P(a_t = a')}{\mu(a' | c_t)} r(c_t, a') \frac{\mathbb{1}\{a' = a\}}{\mu(a' | c_t)} \\ &= r(c_t, a) \end{aligned}$$

$$\hat{V}(\pi) = \frac{1}{T} \sum_{t=1}^T \hat{r}(c_t, \pi(c_t)), \quad \mathbb{E}[\hat{V}(\pi)] = V(\pi)$$

$$\mathbb{E}[(\hat{V}(\pi) - V(\pi))^2] = \frac{1}{T^2} \mathbb{E} \left[\sum_{t=1}^T (\hat{r}(c_t, \pi(c_t)) - r(c_t, \pi(c_t)))^2 \right]$$

$$\leq \frac{1}{T^2} \mathbb{E} \left[\sum_t \hat{r}(c_t, \pi(c_t))^2 \right]$$

$$= \frac{1}{T^2} \sum_t \mathbb{E}_{c_t} \left[r_t^2 \frac{\mathbb{1}\{a_t = \pi(c_t)\}}{P_t^2} \right]$$

$$\leq \frac{1}{T^2} \sum_t \mathbb{E}_{a_t, c_t} \left[\frac{\mathbb{1}\{a_t = \pi(c_t)\}}{\mu(a_t | c_t)^2} \right]$$

$$= \frac{1}{T} \mathbb{E}_{c \sim \nu} \left[\frac{1}{\mu(\pi(c) | c)} \right]$$

I want high prob bound on $|\hat{V}(\pi) - V(\pi)|$.

Hoeffding says that if $Z_t \in [a, b]$ and

$$\mathbb{E}[Z_t] = 0 \quad \text{then} \quad \mathbb{P}\left(\frac{1}{T} \sum_{t=1}^T Z_t \geq |b-a| \sqrt{\frac{\log(1/\delta)}{2T}}\right) \leq \delta.$$

$$\hat{V}(c, \pi(c)) \in \left[0, \frac{1}{\min_{c,a} \mu(a|c)}\right]$$

Bernstein's inequality says that if $Z_t \leq B$ and

$$\mathbb{E}[Z_t] = 0, \quad \mathbb{E}[Z_t^2] \leq \sigma^2 \quad \text{then}$$

$$\mathbb{P}\left(\frac{1}{T} \sum_{t=1}^T Z_t \geq \sqrt{\frac{2\sigma^2 \log(1/\delta)}{T}} + \frac{2B \log(1/\delta)}{3T}\right) \leq \delta.$$

$$\Rightarrow \text{w.p.} \geq 1 - \delta$$

$$|\hat{V}(\pi) - V(\pi)| \leq \underbrace{\sqrt{\mathbb{E}_c \left[\frac{1}{\mu(\pi(c)|c)} \right]}_{\leq \eta \text{ if union}} \cdot \frac{2 \log(2/\delta)}{T} + \underbrace{\max_{c,a} \frac{1}{\mu(a|c)}}_{\leq \eta \text{ if union}} \cdot \frac{2 \log(2/\delta)}{3T}$$

Side note: There exists an estimator $\hat{\mu}: \mathbb{R}^T \rightarrow \mathbb{R}$

s.t. if $\mathbb{E}[Z_t] = 0, \mathbb{E}[Z_t^2] \leq \sigma^2$ then

$$\mathbb{P}\left(\hat{\mu}(\{Z_t\}_{t=1}^T) \geq \sqrt{\frac{2\sigma^2 \log(1/\delta)}{T}}\right) \leq \delta.$$

See Catoni's estimator or median of means.

If $\mu(a|c) = \frac{1}{n}$ for c, a then for a $\pi \in \Pi$

$$|\hat{V}(\pi) - V(\pi)| \leq \sqrt{\frac{2n \log(2/d)}{T}} + \frac{2n \log(2/d)}{3T}$$

$$\leq \sqrt{\frac{4n \log(2/d)}{T}}$$

and $\forall \pi$ w.p. $\geq 1 - \delta$

$$|\hat{V}(\pi) - V(\pi)| \leq \sqrt{\frac{4n \log(2/\pi/d)}{T}} = C$$

You collected a data set computed $\hat{V}(\pi)$

for all $\pi \in \Pi$. Which one do you choose?

Natural to choose $\hat{\pi}_{MLE} = \arg \max_{\pi \in \Pi} \hat{V}(\pi)$

Define $C(\pi) = \sqrt{\mathbb{E}\left[\frac{1}{\mu(\pi(c)|c)}\right] - \frac{2 \log(2/\pi/d)}{T}} + 2 \frac{\max_{c \in \mathcal{C}} \frac{1}{\mu(\pi(c)|c)}}{3T}$

$$V(\hat{\pi}_{MLE}) \geq \hat{V}(\hat{\pi}_{MLE}) - C(\hat{\pi}_{MLE})$$

$$\geq \hat{V}(\pi_{\star}) - C(\hat{\pi}_{MLE})$$

$$\geq V(\pi_x) - C(\pi_x) - C(\hat{\pi})$$

$$\geq V(\pi_x) - 2 \max_{\pi \in \Pi} C(\pi)$$

Pessimism

Define $\hat{\pi}_{\text{pers}} = \underset{\pi \in \Pi}{\text{argmax}} \hat{V}(\pi) - C(\pi)$

$$V(\hat{\pi}_{\text{pers}}) \geq \hat{V}(\hat{\pi}_{\text{pers}}) - C(\hat{\pi}_{\text{pers}})$$

$$\geq \hat{V}(\pi_x) - C(\pi_x)$$

$$\geq V(\pi_x) - 2C(\pi_x)$$

Doubly Robust Estimator

$$\hat{r}_{\text{DR}}(l_t, a) = \hat{f}(l_t, a) + (r_t - \hat{f}(l_t, a)) \frac{\mathbb{1}\{a_t = a\}}{P_t}$$

$$\mathbb{E}[\hat{r}_{\text{DR}}(l_t, a)] = r(l_t, a)$$