

Contextual Bandits

Input: action set \mathcal{X}
 (finite)

for $t=1, 2, \dots$

Nature reveals context $c_t \sim \nu$

Player chooses $x_t \in \mathcal{X}$

Nature reveals reward $y_t \in [0, 1]$, $\mathbb{E}[y_t | c_t, x_t] = r(c_t, x_t)$

Suppose $\text{support}(\nu) =: \mathcal{C}$ and $|\mathcal{C}| < \infty$.

Idea: run MAB on each context $c \in \mathcal{C}$

$$\begin{aligned}\text{Reward} &= \sum_{t=1}^T r(c_t, x_t) \\ &= \sum_{c \in \mathcal{C}} \sum_{t=1}^T \mathbb{I}\{c_t = c\} r(c, x_t)\end{aligned}$$

Regret relative to best action per context

$$\text{Regret} = \sum_{c \in \mathcal{C}} \max_x \sum_{t=1}^T \mathbb{I}\{c_t = c\} (r(c, x) - r(c, x_t))$$

Apply MAB to each. $T_c = \sum_{t=1}^T \mathbb{I}\{c_t = c\}$

$$\leq \sum_{c \in \mathcal{C}} \sqrt{|\mathcal{X}| T_c \log(T_c / \delta)}$$

$$= \sum_{c \in \mathcal{C}} \sqrt{T_c} \cdot \sqrt{|\mathcal{X}| \log(T_c / \delta)}$$

$$\leq \sqrt{\left(\sum_{c \in \mathcal{C}} T_c \right)} \cdot \sum_{c \in \mathcal{C}} |\mathcal{X}| \log(T_c / \delta)$$

Cauchy-Schwarz

$$\leq \sqrt{T |\mathcal{X}| |\mathcal{C}| \log(T / \delta)}$$

$$\text{Reward} \geq \sum_{t=1}^T \max_x r(c_t, x) - \sqrt{T|x|\log(T/\delta)}$$

Another simple idea: ignore the context!

$$\text{Reward} = \sum_{t=1}^T r(c_t, x_t)$$

$$\begin{aligned} \text{Regret} &\leq \max_x \sum_{t=1}^T r(c_t, x) - r(c_t, x_t) \\ &\leq \sqrt{T \cdot |X| \cdot \log(T/\delta)} \end{aligned}$$

$$\Rightarrow \text{Reward} \geq \max_x \sum_{t=1}^T r(c_t, x) - \sqrt{T|x|\log(T/\delta)}$$

Define a policy $\pi: C \rightarrow X$. If value is defined as

$$V(\pi) = \mathbb{E}_{C \sim \nu} [r(c, \pi(c))].$$

Consider a collection of policies Π (assume finite), and at each time t , we choose $\pi_t \in \Pi$ and play $\pi_t(c_t) = x_t$.

$$\text{Policy-Regret} = \max_{\pi \in \Pi} \sum_{t=1}^T r(c_t, \pi(c_t)) - r(c_t, \pi_t(c_t))$$

Ex. For policies that ignore context $\Pi = \{\pi(c) = x : x \in X\}$
 $|\Pi| = |X|$

Ex. For MAB at each context $\Pi = \{\pi(c) = x : x \in X, c \in C\}$
 $|\Pi| = |C| \cdot |X|$

One idea: Treat policy regret as MAB over policies Π .

$$\Rightarrow \text{Regret} \leq \sqrt{|\Pi| \cdot T \log(T)}$$

Suppose we had full information s.t. we observe $y_{t,x}$ w/ $E[y_{t,x}|x, c_t] = r(c_t, x)$ at every time.

With this knowledge, define $\hat{V}(\pi) = \frac{1}{T} \sum_{t=1}^T y_{t,\pi(c_t)}$

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

Apply Hoeffding union bound to have

$$P\left(\bigcup_{\pi \in \Pi} \{ |\hat{V}(\pi) - V(\pi)| \geq \sqrt{\frac{\log(2|\Pi|/\delta)}{2T}} \} \right) \leq \delta.$$

$$P\left(\bigcup_{\pi \in \Pi} \bigcup_{t=1}^T \{ |\hat{V}_t(\pi) - V(\pi)| \geq \sqrt{\frac{\log(2T|\Pi|/\delta)}{2t}} \} \right) \leq \delta$$

Let $\bar{\pi}_t := \arg\max_{\pi \in \Pi} \hat{V}_{t-1}(\pi)$

$$\text{Regret} = \max_{\pi} \left[\sum_{t=1}^T r(c_t, \pi(c_t)) - r(c_t, \bar{\pi}_t(c_t)) \right]$$

$$= \max_{\pi} \sum_{t=1}^T V(\pi) - V(\bar{\pi}_t)$$

$$\stackrel{\bar{\pi}_t = \arg\max_{\pi} V(\pi)}{\leq} 1 + \max_{\pi} \sum_{t=2}^T V(\pi) - V(\bar{\pi}_t)$$

$$= 1 + \sum_{t=2}^T \underbrace{V(\bar{\pi}_t) - \hat{V}_{t-1}(\bar{\pi}_t)}_{\frac{1}{\sqrt{t}}} + \underbrace{\hat{V}_t(\bar{\pi}_t) - \hat{V}_{t-1}(\bar{\pi}_t)}_{\leq 0} + \underbrace{\hat{V}_t(\bar{\pi}_t) - V(\bar{\pi}_t)}_{\leq \frac{1}{\sqrt{t}}}$$

$$\leq 1 + \sum_{t=2}^T 2 \sqrt{\frac{\log(10T/\delta)}{2(t-1)}} \leq \sqrt{T \log(10T/\delta)}.$$

Off-policy evaluation

Fix a logging policy $\mu(\cdot | c) \in \Delta_X$ that prob. dist over actions for each context.

Ex. for any $\lambda \in \Delta_\pi$ I could set $\mu(x|c) = \sum_{\pi \in \Pi} \mathbb{1}\{\pi(c)=x\} \lambda_\pi$

Ex. Take $\mu(x|c) = \frac{1}{|X|}$

Suppose $\mu(\cdot)$ is run for T timesteps to

produce data $\{(c_t, x_t, P_t, y_t)\}_{t=1}^T$ where

$$P_t = \mu(x_t | c_t). \quad x_t \sim \mu(\cdot | c_t)$$

$$\hat{V}_{\text{IPS}}(\pi) := \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{1}\{x_t = \pi(c_t)\}}{P_t} y_t = \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{1}\{x_t = \pi(c_t)\}}{\mu(x_t | c_t)} y_t$$

$$\begin{aligned} \mathbb{E}[\hat{V}_{\text{IPS}}(\alpha)] &= \frac{1}{T} \sum_{t=1}^T \underbrace{\mathbb{E}\left[\frac{\mathbb{1}\{x_t = \pi(c_t)\}}{\mu(x_t | c_t)} y_t\right]} \\ &= \mathbb{E}\left[\underbrace{\mathbb{E}\left[\frac{\mathbb{1}\{x_t = \pi(c_t)\}}{\mu(x_t | c_t)} y_t | x_t, c_t\right]}\right] \\ &= \frac{\mathbb{1}\{x_t = \pi(c_t)\}}{\mu(x_t | c_t)} r(c_t, x_t) \\ &= \mathbb{E}\left[\underbrace{\mathbb{E}\left[\frac{\mathbb{1}\{x_t = \pi(c_t)\}}{\mu(x_t | c_t)} r(c_t, x_t) | c_t = c\right]}\right] \end{aligned}$$

Implicitly assumes $\mu(x|c) \geq 0$ $\forall x, c$

$$\begin{aligned} &\Rightarrow = \sum_{x \in X} \underbrace{\mu(x | c_t)}_{= p(x_t = x | c_t)} \frac{\mathbb{1}\{x = \pi(c_t)\}}{\mu(x | c_t)} r(c_t, x) \\ &= r(c_t, \pi(c_t)) \end{aligned}$$

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

$$\begin{aligned} \mathbb{E}[(x - \mathbb{E}[x])^2] &= \mathbb{E}[x^2] - \mathbb{E}[x]^2 \\ &\leq \mathbb{E}[x^2] &= V(\pi) \end{aligned}$$

$$\text{Consider } \mathbb{V}\left(\hat{V}(\pi)\right) = \frac{1}{T} \text{Var}\left(\underbrace{\frac{\mathbb{I}\{x_t = \pi(c_t)\}}{\mu(x_t | c_t)} y_t}_{\text{y}_t} \right)$$

$$\leq \frac{1}{T} \mathbb{E}\left[\left(\underbrace{\frac{\mathbb{I}\{x_t = \pi(c_t)\}}{\mu(x_t | c_t)} y_t}_{y_t}\right)^2\right]$$

$$y_t \in \{0, 1\} \leq \frac{1}{T} \mathbb{E}\left[\underbrace{\frac{\mathbb{I}\{x_t = \pi(c_t)\}}{\mu(x_t | c_t)^2}}_{\text{y}_t}\right]$$

$$x_t \quad t: 1, 2, \dots, T = \frac{1}{T} \mathbb{E}\left[\sum_x \mu(x | c_t) \frac{\mathbb{I}\{x | \pi(c_t)\}}{\mu(x | c_t)^2}\right]$$

$$\mathbb{V}\left(\frac{1}{T} \sum_t x_t\right) = \frac{1}{T^2} \mathbb{V}\left(\sum_t x_t\right)$$

$$= \frac{1}{T^2} \sum_t \mathbb{V}(x_t) = \frac{1}{T} \mathbb{E}\left[\sum_x \frac{\mathbb{I}\{x | \pi(c_t)\}}{\mu(x | c_t)}\right]$$

$$= \frac{1}{T} \mathbb{V}(x_t) = \frac{1}{T} \mathbb{E}\left[\frac{1}{\mu(\pi(c_t) | c_t)}\right]$$

for any t .

Lemma | Bernstein's inequality. Let X_1, \dots, X_n be independent R.V. s.t. $\frac{1}{n} \sum_{i=1}^n \mathbb{E}[(X_i - \mathbb{E}[X_i])^2] \leq \sigma^2$ and $\max_i |X_i| \leq B$. Then

$$\text{w.p. } \geq 1 - \delta : \left| \frac{1}{n} \sum_{i=1}^n (X_i - \mathbb{E}[X_i]) \right| \leq \sqrt{\frac{2\sigma^2 \log(2/\delta)}{n}} + \frac{2B \log(2/\delta)}{3n}$$

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

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

Ex. $\mu(x|c) = \frac{1}{|x|}$ then w.p. $\geq 1-\delta$ if $c \in \Pi$

$$\begin{aligned} |\hat{V}_{IPS}(\pi) - V(\pi)| &\leq \sqrt{\frac{2|x|\log(2|\Pi|/\delta)}{T}} + \frac{2|x|\log(2/\delta)}{3T} \\ &\leq \sqrt{\frac{4|x|\log(2|\Pi|/\delta)}{T}} \end{aligned}$$

Above is "model-free" approach to estimating $V(\pi)$

in the sense that it makes no attempt to "model the environment" meaning estimating $r(c, x) \forall c, x$.

Example policy class given $\Phi : \mathcal{C} \times \mathcal{X} \rightarrow \mathbb{R}^d$, for each $\theta \in \mathbb{R}^d$

Define $T_\theta(c) = \arg\max_x \langle \phi(c, x), \theta \rangle$

Fact after T rounds $|\Pi| \leq n^d$ policies

Model-Based estimation

If you know $r(c, x)$ then best possible action (and policy) is to play $x_t = \arg\max_x r(c_t, x)$.

Next best thing: consider a function class \mathcal{F} so that $f \in \mathcal{F}$ has $f: \mathcal{C} \times \mathcal{X} \rightarrow \mathbb{R}$, given data $\{(c_t, x_t, p_t, y_t)\}_{t=1}^T$

$$\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{t=1}^T (f(c_t, x_t) - y_t)^2$$

set $\hat{V}_{\mathcal{F}}(\pi) = \frac{1}{T} \sum_{t=1}^T \hat{f}(c_t, \pi(c_t))$.

Ex. Let $\mathcal{F} = \{\langle \phi(c, x), \theta \rangle : \theta \in \mathbb{R}^d\}$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum_{t=1}^T (\langle \phi(c_t, x_t), \theta \rangle - y_t)^2$$

$$\hat{V}_{\mathcal{F}}(\pi) = \frac{1}{T} \sum_{t=1}^T \langle \phi(c_t, \pi(c_t)), \hat{\theta} \rangle$$