

# Policy Bandit Optimization

for  $t=1, 2, \dots$

Nature reveals  $c_t$

Player chooses  $x_t$

Nature reveals  $y_t$  s.t.  $\mathbb{E}[y_t | c_t, x_t] = r(c_t, x_t)$

Given policy set  $\Pi$ , minimize policy regret

$$\max_{\pi} \sum_{t=1}^T r(c_t, \pi(c_t)) - r(c_t, x_t)$$

$$V_{\pi} = \operatorname{argmax}_{\pi} V(\pi)$$

## $\mathfrak{S}$ -greedy

for  $t=1, 2, \dots, \mathfrak{S}$  play  $x_t \sim \text{uniform}(X)$  and then set

$$\hat{\pi} = \operatorname{argmax}_{\pi} \hat{V}_{\text{IPS}}(\pi) = \operatorname{argmax}_{\pi} \frac{1}{\mathfrak{S}} \sum_{t=1}^{\mathfrak{S}} \frac{\mathbb{1}\{x_t = \pi(c_t)\}}{p_t} y_t$$

$$p_t = \frac{1}{|X|}$$

By Bernstein's inequality, w.p.  $\geq 1 - \delta$

$$|\hat{V}_{\text{IPS}}(\pi) - V(\pi)| \leq \sqrt{\frac{2 \mathbb{E} \left[ \frac{1}{\mu(\pi(c) | c)} \right] \log(|\Pi|/\delta)}{\mathfrak{S}}} + \frac{2 \log(|\Pi|/\delta)}{\mathfrak{S} \cdot \underbrace{\min_{x,c} \mu(x|c)}_{= \frac{1}{|X|}}}$$

$= |X|$

$$\leq c \sqrt{\frac{|X| \log(|\Pi|/\delta)}{\mathfrak{S}}}$$

$$\text{Regret} = \sum_{t=1}^T V(\pi_{\hat{\pi}}) - r(c_t, x_t)$$

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

$$= \left( \sum_{t=1}^3 V(\pi_a) - r(c_t, x_t) \right) + \underbrace{\sum_{t=3+1}^T V(\pi_a) - r(c_t, \hat{\pi}(c_t))}_{(T-3) (V(\pi_a) - V(\hat{\pi}))}$$

$$\leq 3 + (T-3) (V(\pi_a) - V(\hat{\pi}))$$

$$V(\pi_a) - V(\hat{\pi}) = V(\pi_a) - \hat{V}(\pi_a) + \underbrace{\hat{V}(\pi_a) - \hat{V}(\hat{\pi})}_{\leq 0} + \hat{V}(\hat{\pi}) - V(\hat{\pi})$$

$$\leq 2c \sqrt{\frac{|\chi| \log(2/\delta)}{3}}$$

$$\rightarrow \leq 3 + T \cdot 2c \sqrt{\frac{|\chi| \log(2/\delta)}{3}}$$

Set  $\mathfrak{J} = (|\chi| T^2 \log(1/\delta))^{1/3}$  then

$$\text{Regret} \leq T^{2/3} (|\chi| \log(1/\delta))^{1/3}$$

But how do you find  $\hat{\pi}$ ?

$$\begin{aligned} \hat{\pi} &= \operatorname{argmax}_{\pi} \hat{V}_{\text{IPS}}(\pi) = \operatorname{argmax}_{\pi} \frac{1}{3} \sum_{t=1}^3 \frac{\mathbb{1}\{x_t = \pi(c_t)\}}{P_t} y_t \\ &= \operatorname{argmax}_{\pi} \frac{1}{3} \sum_{t=1}^3 \frac{\mathbb{1}\{x_t = \pi(c_t)\}}{P_t} y_t \\ &= \operatorname{argmax}_{\pi} \frac{1}{3} \sum_{t=1}^3 \frac{(1 - \mathbb{1}\{x_t \neq \pi(c_t)\})}{P_t} y_t \\ &= \operatorname{argmax}_{\pi} -\frac{1}{3} \sum_{t=1}^3 \frac{\mathbb{1}\{x_t \neq \pi(c_t)\}}{P_t} y_t \end{aligned}$$

$$= \operatorname{argmin}_{\pi} \frac{1}{3} \sum_{t=1}^T \frac{y_t}{P_t} \mathbb{1}\{\pi(c_t) \neq x_t\}$$

Problem: a minimizer of this is  $\pi(c_t) \neq x_t$  for all  $c_t$ . Not necessarily converging to  $\pi_*$ !

- Regularization
- Smaller TTT
- Doubly robust estimators.

## Elimination alg for $\sqrt{T}$ -regret bound

Input  $\Pi$

Init  $\Pi_1 = \Pi, T_1 = 0$

for  $l=1, 2, \dots$

$\varepsilon_l = 2^{-l}$ . Set  $\beta_l = 16 \varepsilon_l^{-2} / \chi \log\left(\frac{4R^2 |\Pi|}{\delta}\right)$ ,  $T_l = T_{l-1} + \beta_l$

Define  $\lambda \in \Delta_{\Pi_l}$ ,  $Q(x|c) = \sum_{\pi \in \Pi} \lambda_{\pi} \mathbb{1}\{\pi(c) = x\}$

for  $t = T_{l-1} + 1, \dots, T_l$  where  $\gamma = \min\left\{\frac{1}{2|\chi|}, \sqrt{\frac{\log(2/(\delta|\chi|))}{4|\chi|\beta_l}}\right\}$

Nature reveals context  $c_t$

Draw  $\tilde{\pi}_t \sim \lambda$ , Play  $x_t = \begin{cases} \tilde{\pi}_t(c_t) & \text{w.p. } 1 - \delta|\chi| \\ \text{unif}(\chi) & \text{w.p. } \delta|\chi| \end{cases}$   $P_t = \gamma + (1 - \delta|\chi|)Q(x_t|c_t)$

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

$$\hat{V}(\pi) = \frac{1}{T_l - T_{l-1}} \sum_{t \in T_l} y_t \frac{\mathbb{1}\{\pi(c_t) = x_t\}}{P_t}$$

$$\Pi_{l+1} = \Pi_l \setminus \left\{ \pi \in \Pi_l : \max_{\pi'} \hat{V}(\pi') - \hat{V}(\pi) > 2\varepsilon_l \right\}$$

Define  $Q(x|c) = \sum_{\pi \in \Pi} \lambda_{\pi} \mathbb{1}\{\pi(c) = x\}$

$$\begin{aligned} & \min_{\lambda \in \Delta_{\Pi}} \max_{\pi \in \Pi} \mathbb{E} \left[ (\hat{V}(\pi) - V(\pi))^2 \right] \\ &= \min_{\lambda \in \Delta_{\Pi}} \max_{\pi \in \Pi} \mathbb{E}_c \left[ \frac{1}{Q(\pi(c)|c)} \right] \end{aligned}$$

Lemma For any policy set  $\Pi$ ,

$$\min_{\lambda \in \Delta_{\Pi}} \max_{\pi \in \Pi} \mathbb{E} \left[ \frac{1}{Q(\pi(c)|c)} \right] \leq |\mathcal{X}|.$$

Moreover, if  $Q^{\gamma}(x|c) := \gamma + (1 - \gamma|\mathcal{X}|) Q(x|c)$  then

$$Q^{\gamma}(x|c) \geq \gamma \quad \text{but also for } \gamma \leq \frac{1}{2|\mathcal{X}|}$$

$$\min_{\lambda \in \Delta_{\Pi}} \max_{\pi \in \Pi} \mathbb{E} \left[ \frac{1}{Q^{\gamma}(\pi(c)|c)} \right] \leq 2|\mathcal{X}|.$$

---


$$|\hat{V}^{\pi} - V(\pi)| \leq \sqrt{\frac{\mathbb{E} \left[ \frac{1}{Q^{\gamma}(\pi(c)|c)} \right] 2 \log(2|\Pi|/\delta)}{\mathfrak{I}_L}} + \frac{2 \log(2|\Pi|/\delta)}{3 \min_{c, x} Q^{\gamma}(x|c) \mathfrak{I}_L}$$

$$\leq \sqrt{\frac{4|\mathcal{X}| \log(2|\Pi|/\delta)}{\mathfrak{I}_L}} + \frac{2 \log(2|\Pi|/\delta)}{3\gamma \mathfrak{I}_L}$$

$$\leq \sqrt{\frac{16|\mathcal{X}| \log(2|\Pi|/\delta)}{\mathfrak{I}_L}} \quad \text{if } \gamma = \min \left\{ \frac{1}{2|\mathcal{X}|}, \sqrt{\frac{\log(2|\Pi|/\delta)}{4|\mathcal{X}| \mathfrak{I}_L}} \right\}$$

$$\leq \epsilon_e \quad \text{if} \quad \mathcal{J}_e \geq 16 \epsilon_e^{-2} |\chi| \log(2|\mathcal{T}|/\delta)$$

Lemma w.p.  $\geq 1-\delta$ ,  $\forall e \pi_e \in \Pi_e$  and

$$\max_{\bar{a} \in \Pi_e} V(\bar{a}_e) - U(\pi_e) \leq 8\epsilon_e.$$

$$\text{Regret} = \sum_{t=1}^T V(\pi_{e_t}) - Y_t$$

$$\leq \sum_{e=1}^L \mathcal{J}_e \left( \gamma |\chi| + \overset{\leq 1/2}{(1-\gamma)\delta} 8\epsilon_e \right)$$

$$\leq \sum_{e=1}^L \mathcal{J}_e \left( |\chi| \cdot \sqrt{\frac{\log(1|\mathcal{T}|l^2/\delta)}{|\chi| \mathcal{J}_e}} + 4\epsilon_e \right)$$

$$= \sum_{e=1}^L \mathcal{J}_e \left( c \epsilon_e + 4\epsilon_e \right)$$

$\downarrow$  Plug in  $\mathcal{J}_e = c \epsilon_e^{-2} |\chi| \log(1|\mathcal{T}|l^2/\delta)$

$$\leq \sum_{e=1}^L \min \left\{ \epsilon_e \mathcal{J}_e, \epsilon_e \cdot c \epsilon_e^{-2} |\chi| \log(1|\mathcal{T}|l^2/\delta) \right\}$$

$$\leq \sum_{e=1}^L \min \left\{ \epsilon_e \mathcal{J}_e, c \epsilon_e^{-1} |\chi| \log(1|\mathcal{T}|l^2/\delta) \right\}$$

$$\leq \sum_{e=1}^L \max_{\epsilon} \min \left\{ \epsilon \mathcal{J}_e, c \epsilon^{-1} |\chi| \log(1|\mathcal{T}|l^2/\delta) \right\}$$

$$\langle x, y \rangle \leq \|x\| \cdot \|y\|$$

$$= \sum_{t=1}^L \sqrt{c} \sqrt{|\mathcal{X}| \log_2 \left( \frac{1}{\pi} \frac{L^2}{\delta} \right)} \cdot \sqrt{\beta_{\epsilon}}$$

$$\leq \sqrt{L \cdot c |\mathcal{X}| \log_2 \left( \frac{1}{\pi} \frac{L^2}{\delta} \right)} \cdot \underbrace{\sqrt{\sum_{t=1}^L \beta_{\epsilon}}}_{= \sqrt{T}}$$

$$\leq \sqrt{c L |\mathcal{X}| T \log_2 \left( \frac{1}{\pi} \frac{L^2}{\delta} \right)} \quad L = \log_2 \left( \frac{1}{\Delta_{\min}}, T \right)$$

## Model-based Contextual Bandits

Given  $\mathcal{F}$  and data  $\{(c_t, x_t, p_t, y_t)\}$  fit

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

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

Ex. Linear policy class:  $f(c, x) = \phi(x, c)^T \theta$  for  $\theta \in \mathbb{R}^d$

and  $\phi: \mathcal{C} \times \mathcal{X} \rightarrow \mathbb{R}^d$  is a feature map.

Ex. suppose  $C_t \in \mathbb{R}^p$  and  $\mathcal{X} = \{1, \dots, n\}$  then

natural choice  $\phi(c, x) = \text{vec}(c e_x^\top) \in \mathbb{R}^{pn}$

$$= \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & c \end{bmatrix}$$

Algorithm: UCB  $\hat{\theta} = \underset{\theta}{\text{argmin}} \sum_{t=1}^T (y_t - \phi(c_t, x_t)^\top \theta)^2 + \gamma \|\theta\|^2$

**Proposition 5.** Fix  $\delta \in (0, 1)$ ,  $\gamma \geq 0$ , and  $\theta_* \in \mathbb{R}^d$ . Assume for all  $t \geq 1$  that  $y_t = \langle \theta_*, x_t \rangle + \eta_t$  and  $\mathbb{E}[\exp(s\eta_t) | \mathcal{F}_{t-1}] \leq \exp(s^2/2)$  for any  $s \in \mathbb{R}$  where  $\mathcal{F}_t$  is such that  $x_1, y_1, \dots, x_{t-1}, y_{t-1}, x_t$  are  $\mathcal{F}_{t-1}$  measurable. If  $S_t = \sum_{i=1}^t x_i \eta_i$ ,  $V_t = \sum_{i=1}^t x_i x_i^\top$ , and  $\hat{\theta}_t = (V_t + \gamma I)^{-1} S_t$ , then

$$\|\hat{\theta}_t - \theta_*\|_{(V_t + \gamma I)} \leq \sqrt{\gamma} \|\theta_*\|_2 + \sqrt{2 \log(1/\delta) + \log(\gamma^{-d} |V_t + \gamma I|)}$$

for all  $t \geq 1$  simultaneously with probability at least  $1 - \delta$ . Moreover, if  $\max_t \|x_t\|_2^2 \leq L$  then  $\log(\gamma^{-d} |V_t + \gamma I|) \leq d \log(\frac{tL}{d\gamma} + 1)$ .

for  $t = 1, 2, \dots$

Nature reveals context  $c_t$

Player chooses  $x_t = \underset{x \in \mathcal{X}}{\text{argmax}} \text{UCB}(c_t, x)$

where  $\text{UCB}(c_t, x) = \max_{\theta \in \mathcal{E}_t} \langle \theta, \phi(c_t, x) \rangle$

$$\mathcal{E}_t = \left\{ \theta : \|\hat{\theta}_t - \theta\|_{V_t + \gamma I} \leq \dots \right\}$$

Observe  $y_t$ ,

Update  $\hat{\theta}_t, V_t$

# Thompson Sampling

for  $t = 1, 2, 3, \dots$

Nature reveals  $c_t$

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

Play  $x_t = \underset{x \in \mathcal{X}}{\operatorname{argmax}} \phi(c_t, x)^T \theta'$