# Policy Bandit Optimization

for $t = 1, 2, \ldots$

    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

$$\max_{\pi} \sum_{t=1}^{T} r(c_t, \pi(c_t)) - r(c_t, x_t) \qquad \pi_\star = \text{argmax}_{\pi} V(\pi)$$

## $\varsigma$-greedy

for $t = 1, 2, \ldots, \varsigma$   play $x_t \sim \text{uniform}(x)$ and then set

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

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

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

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

$$\underbrace{= |x|}_{} \qquad \underbrace{= \frac{1}{|x|}}_{}$$

$$\leq_c \sqrt{\frac{|x| \log(2|\Pi|/\delta)}{\varsigma}}$$

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

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

$$= \left( \sum_{t=1}^{\mathcal{I}} V(\pi_\Delta) - r(c_t, x_t) \right) + \sum_{t=\mathcal{I}+1}^{T} V(a_\ge) - r(c_t, \hat{\pi}(c_t))$$

$$\le \quad \mathcal{I} \quad + \quad \underbrace{(T-\mathcal{I})\left( V(\pi_\Delta) - V(\hat{\pi}) \right)}$$

$$V(\pi_\Delta) - V(\hat{\pi}) = V(\pi_\Delta) - \hat{V}(\pi_\Delta) + \underbrace{\hat{V}(\pi_\Delta) - \hat{V}(\hat{\pi})}_{\color{red}{\le 0}} + \hat{V}(\hat{\pi}) - V(\hat{\pi})$$

$$\le 2c \sqrt{\frac{|\mathcal{X}| \log(2|\pi|/\delta)}{\mathcal{I}}}$$

$$\le \quad \mathcal{I} + T \cdot 2c \sqrt{\frac{|\mathcal{X}| \log(2|\pi|/\delta)}{\mathcal{I}}}$$

Set $\quad \mathcal{I} = \left( |\mathcal{X}| T^2 \log(|\pi|/\delta) \right)^{1/3} \quad$ then

$$\text{Regret} \le T^{2/3} \left( |\mathcal{X}| \log(|\pi|/\delta) \right)^{1/3}$$

___

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

$$\hat{\pi} = \operatorname*{argmax}_\pi \hat{V}_{IPS}(\pi) = \operatorname*{argmax}_\pi \frac{1}{\mathcal{I}} \sum_{t=1}^{\mathcal{I}} \frac{\mathbb{1}\{x_t = \pi(c_t)\}}{P_t} y_t$$

$$= \operatorname*{argmax}_\pi \frac{1}{\mathcal{I}} \sum_{t=1}^{\mathcal{I}} \frac{\mathbb{1}\{x_t = \pi(c_t)\}}{P_t} y_t$$

$$= \operatorname*{argmax}_\pi \frac{1}{\mathcal{I}} \sum_{t=1}^{\mathcal{I}} \frac{(1 - \mathbb{1}\{x_t \ne \pi(c_t)\})}{P_t} y_t$$

$$= \operatorname*{argmax}_\pi -\frac{1}{\mathcal{I}} \sum_{t=1}^{\mathcal{I}} \frac{\mathbb{1}\{x_t \ne \pi(c_t)\}}{P_t} y_t$$

$$= \underset{\pi}{\text{argmin}} \ \frac{1}{J} \sum_{t=1}^{J} \frac{y_t}{p_t} \mathbb{1}\{x_t \neq \pi(c_t)\}$$

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

- Regularization
- Smaller $|\Pi|$
- Doubly robust estimators.

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

Input $\Pi$
Init $\Pi_1 = \Pi, \ T_1 = 0$
for $\ell = 1, 2, \ldots$

$\quad \varepsilon_\ell = 2^{-\ell}$. Set $S_\ell = \ldots$ $\qquad\qquad , \ T_\ell = T_{\ell-1} + S_\ell$

$\quad$ Define $\lambda \in \Delta_{\Pi_\ell}$

$\quad$ for $t = T_{\ell-1} + 1, \ldots, T_\ell$

$\qquad$ Nature reveals context $c_t$

$\qquad$ Draw $\pi_t \sim \lambda$, play $x_t = \pi(c_t)$ set $p_t = \sum_{\pi \in \Pi} \lambda_\pi \mathbb{1}\{x_t = \pi(c_t)\}$

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

$$\hat{V}(\pi) = \frac{1}{T_\ell - T_{\ell-1}} \sum_{t=T_{\ell-1}}^{T_\ell} y_t \frac{\mathbb{1}\{\pi(c_t) = x_t\}}{p_t}$$

$$\Pi_{\ell+1} = \Pi_\ell \setminus \{\pi \in \Pi_\ell : \underset{\pi'}{\max} \ \hat{V}(\pi') - \hat{V}(\pi) > 2\varepsilon_\ell\}$$

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

$$\min_{\lambda \in \Delta_\Pi} \max_{\pi \in \Pi} \mathbb{E}\left[ \left( \hat{V}(\pi) - V(\pi) \right)^2 \right]$$

$$= \min_{\lambda \in \Delta_\Pi} \max_{\pi \in \Pi} \mathbb{E}_c\left[ \frac{1}{Q(\pi(c) \mid c)} \right]$$

**Lemma** / For any policy set $\Pi$,

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