

A \sqrt{T} regret algorithm for contextual bandits

By defining $\mu(a|c)$ we can estimate $\hat{V}(\pi)$ like

$$\text{Variance}(\hat{V}(\pi)) = \mathbb{E} \left[\frac{1}{\mu(a|c)} \right]$$

Elimination alg.

Input: $\Pi, \delta \in (0,1)$

Init: $\Pi_1 = \Pi, \ell = 1, T_0 = 0$

while $|\Pi_\ell| > 1$

$$\varepsilon_\ell = \frac{\delta}{2^\ell}, \quad \mathcal{I}_\ell = \lceil 32|\mathcal{A}| \varepsilon_\ell^{-2} \log(2|\Pi|/\delta) \rceil, \quad \gamma_\ell = \min \left\{ \frac{1}{2|\mathcal{A}|}, \sqrt{\frac{\log(2|\Pi|/\delta)}{4|\mathcal{A}| \mathcal{I}_\ell}} \right\},$$

$$\mu_\ell(a|c) = \gamma_\ell + (1 - \gamma_\ell |\mathcal{A}|) \sum_{\pi \in \Pi_\ell} \lambda_\pi^{\ell} \mathbb{1}\{\pi(c) = a\}$$

$$\text{where } \lambda^{\ell} = \underset{\lambda \in \Delta_{\Pi_\ell}}{\text{argmin}} \max_{\pi \in \Pi_\ell} \mathbb{E}_c \left[\frac{1}{\sum_{\pi' \in \Pi_\ell} \mathbb{1}\{\pi'(c) = \pi(c)\}} \right]$$

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

Nature reveals context c_t

Algorithm plays $a_t \sim \mu_\ell(\cdot | c_t)$, set $P_t = \mu_\ell(a_t | c_t)$

and obtains r_t w/ $\mathbb{E}[r_t | c_t, a_t] = v(c_t, a_t)$

$$\hat{V}_\ell(\pi) = \frac{1}{T_\ell - T_{\ell-1}} \sum_{t=T_{\ell-1}+1}^{T_\ell} \frac{\mathbb{1}\{\pi(c_t) = a_t\}}{P_t} r_t \quad \forall \pi$$

$$\Pi_{\ell+1} = \Pi_\ell \setminus \left\{ \pi \in \Pi_\ell : \max_{\pi' \in \Pi_\ell} \hat{V}_\ell(\pi') - \hat{V}_\ell(\pi) > 2\varepsilon_\ell \right\}$$

$\ell \leftarrow \ell + 1$

Want to choose $\mu_\ell(a|c)$ and \mathcal{I}_ℓ s.t.

$$|\hat{V}_\ell(\pi) - V(\pi)| \leq \varepsilon_\ell \quad \forall \pi \quad (\text{with high prob})$$

and incur $\leq \varepsilon_\ell$ regret per time step t .

What if $\mu_c(a|c) \propto \mathbb{1}\{\exists \pi \in \Pi_c: \pi(c)=a\}$

Very reasonable for exploration but we cannot (easily) bound the regret of the policy, even though we know (or will show) $|\hat{V}_t(\pi) - V(\pi)| \leq \delta \epsilon_c$.

Alternatively, If $|\hat{V}_t(\pi) - V(\pi)| \leq \epsilon_c$ then by elim rule

$$\max_{\pi \in \Pi_c} V(\pi_*) - V(\pi) \leq \delta \epsilon_c.$$

If I choose

any $\pi_t \in \Pi_c$ ^{randomly} according to some distribution $\lambda \in \Delta_{\Pi_c}$

then $\mathbb{E}_{\pi_t \sim \lambda} [V(\pi_*) - V(\pi_t)]$

$$= \sum_{\pi \in \Pi_c} \lambda_\pi (V(\pi_*) - V(\pi))$$

$$\leq \sum_{\pi \in \Pi_c} \lambda_\pi \delta \epsilon_c = \delta \epsilon_c$$

Idea: draw $\pi_t \sim \lambda_\pi$ and play $a_t = \pi_t(c_t)$

$$\mu_c(a|c) = \sum_{\pi \in \Pi_c} \lambda_\pi \mathbb{1}\{\pi(c)=a\}$$

How do we choose λ ? Well we want to minimize \mathcal{I}_c : estimate each π ϵ_c accuracy.

If $\gamma_e = \min_{a,c} \mu_e(a|c)$ then

$$|\hat{V}_e(\pi) - V(\pi)| \leq \sqrt{\mathbb{E}_c \left[\frac{1}{\mu_e(\pi(c)|c)} \right] \frac{2 \log(2/\pi/d)}{3_e}} + \frac{2 \log(2/\pi/d)}{2 3_e \gamma_e}$$

$$\leq \sqrt{\mathbb{E}_c \left[\frac{1}{\mu_e(\pi(c)|c)} \right] \frac{16 \log(2/\pi/d)}{3_e}} \quad \begin{array}{l} \text{Want w/ minimum } 3_e \\ \leq \epsilon_e \end{array}$$

Holds for conditions on γ_e .

Choose $\lambda \in \Delta_{\pi_e}$ to minimize where $\mu_e(a|c) = \sum_{\pi} \lambda_{\pi} \mathbb{1}\{\pi(c)=a\}$

Claim for any set of policies $\Pi : \pi : C \rightarrow A$, $\gamma \leq \frac{1}{2|A|}$

$$\min_{\lambda \in \Delta_{\Pi}} \max_{\pi \in \Pi} \mathbb{E}_c \left[\frac{1}{\sum_{\pi'} \lambda_{\pi'} \mathbb{1}\{\pi'(c)=\pi(c)\}} \right] \leq |A|$$

Moreover if $\mu(a|c) = \gamma + (1 - |A|\gamma) \sum_{\pi'} \lambda_{\pi'} \mathbb{1}\{\pi'(c)=a\}$

$$\max_{\pi \in \Pi} \mathbb{E}_c \left[\frac{1}{\mu(a|c)} \right] \leq \max_{\pi \in \Pi} \mathbb{E}_c \left[\frac{2}{\sum_{\pi'} \lambda_{\pi'} \mathbb{1}\{\pi'(c)=\pi(c)\}} \right] \leq 2|A|$$

Suppose it is. Then

$$\max_{\pi} |\hat{V}_e(\pi) - V(\pi)| \leq \max_{\pi} \sqrt{\mathbb{E}_c \left[\frac{1}{\mu_e(\pi(c)|c)} \right] \frac{2 \log(2/\pi/d)}{3_e}} + \frac{2 \log(2/\pi/d)}{2 3_e \gamma_e}$$

For our choice of γ_e →

$$\leq \max_{\pi} \sqrt{\max \left\{ \mathbb{E}_c \left[\frac{1}{\mu_e(\pi(c)|c)} \right], |A| \right\} \frac{16 \log(2/\pi/d)}{3_e}}$$

$$\leq \sqrt{|A| \frac{32 \log(2/\pi/d)}{3_e}}$$

$\leq \epsilon_l$ by our choice of \mathcal{I}_l .

Assume this holds for all l

1) Show $\pi_{\#} \in \Pi_l$ for all l

2) Show $\max_{\pi \in \Pi_l} V(\pi_{\#}) - V(\pi) \leq 8\epsilon_l$

3) Prove regret bound

$$\text{Regret } \nu T + \sum_{l=1}^{\lfloor L \log_2 \frac{1}{\nu} \rfloor} \mathcal{I}_l \left(|A| \gamma_l \cdot 1 + (1 - |A| \gamma_l) \cdot 8\epsilon_l \right)$$

$$\leq \nu T + \sum_{l=1}^{\lfloor L \log_2 \frac{1}{\nu} \rfloor} \underbrace{\mathcal{I}_l |A| \gamma_l}_{\text{Way these are chosen}} + 8\epsilon_l \mathcal{I}_l$$

$$\leq \nu T + C \sum_{l=1}^{\lfloor L \log_2 \frac{1}{\nu} \rfloor} \epsilon_l \mathcal{I}_l$$

$$\leq \nu T + C \log(|\mathcal{M}|/\delta) \sum_{l=1}^{\lfloor L \log_2 \frac{1}{\nu} \rfloor} 2^l$$

$$\leq \nu T + |\mathcal{M}| \log(|\mathcal{M}|/\delta) / \nu$$

$$\text{Optimize over } \nu \Rightarrow \sqrt{|\mathcal{M}| T \log(|\mathcal{M}|/\delta)}$$