

CB continued

Last time -

Explaining method called Square CB

'Reduce''

Inverse Proportional

Contextual
Bandits

Gap weighting

Online
Regression

(we have
studied)

ϵ -Greedy Strategy for CB.

$\epsilon \in [0, 1]$

For $t=1$ to T_0

- Nature gives a context x_t .
- Use "oracle" to forecast rewards for all arms.

$\hat{f}_t : [K] \rightarrow [0, 1]$ estimate (based on x_t & past)

$$f^*(\pi_t, \pi_t) = \mathbb{E}[r_t | x_t, \pi_t] \text{ UNKNOWN}$$

- $\pi_t \sim p_t$
- With probability ϵ , pick $\pi_t \sim \text{Uni}[K]$
 - $1-\epsilon$, pick $\pi_t = \underset{\pi}{\text{argmax}} \hat{f}_t(\pi)$.
 - Observe r_t of our action.

$\text{Pr} \text{im:}$ given an oracle achieving

$$\sum_{t=1}^T \mathbb{E} \left[\left(\hat{f}_t(x_t, \pi_t) - f^*(x_t, \pi_t) \right)^2 \right] \leq \text{Reg}_{SQ}$$
 (typically $\mathcal{O}(T)$)

Then ϵ -Greedy achieves (for some ϵ)

$$\mathbb{E}[\text{Reg}_{CB}] \leq K^{1/3} T^{2/3} \text{Reg}_{SQ}^{1/3}$$

$$\sum_{t=1}^T \left[f^*(x_t, \pi_t^*) - f^*(x_t, \pi_t) \right]$$

where $\pi_t^* = \underset{\pi}{\text{argmax}} f^*(x_t, \pi)$

PF:

$$\mathbb{E} \sum_{t=1}^T [f^*(x_t, \pi_t^*) - f^*(x_t, \pi_t)]$$

$$= \mathbb{E} \sum_{t=1}^T \left[\begin{aligned} & \textcircled{I} \left(\hat{f}(x_t, \pi_t^*) - \hat{f}_t(x_t, \pi_t) \right) \\ & + \left(f^*(x_t, \pi_t^*) - \hat{f}_t(x_t, \pi_t^*) \right) \\ & + \left(\hat{f}_t(x_t, \pi_t) - f^*(x_t, \pi_t) \right) \end{aligned} \right]$$

Observe: $\textcircled{I} = \sum (\hat{f}_t(\pi_t^*) - \hat{f}_t(\pi_t))$
 $\leq \varepsilon T \quad \text{I} = \varepsilon T$

$$|\langle u, v \rangle| \leq \|u\|_2 \|v\|_2 \text{ c.s.}$$

(II)

$$\sum_{t=1}^T (\hat{f}_t(\pi_t) - f^*(\pi_t))$$

$$\leq \sqrt{T \sum_{t=1}^T \mathbb{E} \left[(\hat{f}_t(\pi_t) - f^*(\pi_t))^2 \right]}$$

$$\leq \sqrt{T} \int \text{RegSQ}$$

II

$$\sum_{t=1}^T \left(f^*(x_t, \pi_t^*) - \hat{f}_t(\pi_t^*) \right)$$

$$= \sum_{t=1}^T \frac{\sqrt{p_t(\pi_t^*)}}{\sqrt{p_t(\pi_t^*)}} \left(f^*(x_t, \pi_t^*) - \hat{f}_t(\pi_t^*) \right)$$

$$p_t(\pi_t^*) \geq \frac{\epsilon}{K}$$

$$\leq \sqrt{\sum_{t=1}^T \frac{1}{p_t(\pi_t^*)}} \sqrt{\sum_{t=1}^T p_t(\pi_t^*) \left(f^*(x_t, \pi_t^*) - \hat{f}_t(\pi_t^*) \right)^2}$$

$$\leq \sqrt{\frac{KT}{\epsilon}}$$

$$\leq \sqrt{\text{Reg}_Q}$$

$$\textcircled{I} + \textcircled{II} + \textcircled{III} \leq \varepsilon T + \sqrt{\frac{KT}{\varepsilon}} \sqrt{K_{\text{reg}} \sigma} \left(+ \sqrt{T} \sqrt{K_{\text{reg}} \sigma} \right)$$

$$\text{Let } \varepsilon = \frac{K^{1/3} K_{\text{reg}}^{1/3} \sigma^{1/3}}{T^{1/3}}$$

$$\rightsquigarrow \leq 3 K^{1/3} K_{\text{reg}}^{1/3} \sigma^{1/3} T^{2/3}$$

Alg Square CB:

$\delta > 0$ (parameter to be optimized later)

For $t=1$ to T :

- Nature reveals \hat{x}_t

- Oracle forecast is $f_t: [K] \rightarrow [0, 1]$

- Select $\lambda \in [1, K]$ s.t.

$$\sum_{\pi \in [K]} p_t(\pi) = 1$$

$$p_t(\pi) = \frac{1}{\lambda} \exp(\hat{f}_t(\hat{\pi}_t) - f_t(\pi))$$

- Play $\pi_t \sim p_t$, get r_t .

$$\hat{\pi}_t = \arg \max_{\pi} \hat{F}(\pi)$$

Thm As before, assume oracle

$$\mathbb{E} \left(f_{\tau}(x_{\tau}) - f^*(x_{\tau}, \bar{x}_{\tau}) \right)^2 \leq \text{Reg}_{\text{SQ}}$$

For some $\delta > 0$,

$$\mathbb{E} \text{Reg}_{\text{CB}} \leq \sqrt{KT \text{Reg}_{\text{SQ}}}$$

pf?

Use

$$|a-b| \leq \frac{a^2}{2} + \frac{b^2}{2}$$

$$\leftrightarrow (a-b)^2 \geq 0$$

(AM-GM
Inequality)

in a clever way.

Lemma:

$$\hat{f}: [K] \rightarrow [0, 1]$$

$$f^*: [K] \rightarrow [0, 1]$$

$$\lambda \in [1, A]$$

$$\hat{\pi} = \operatorname{argmax}_{\pi} \hat{f}(\pi)$$

$$\pi^* = \operatorname{argmax}_{\pi} f^*(\pi)$$

$$p(\pi) = \frac{1}{\lambda + 2\sigma(\hat{f}(\pi) - f^*(\pi))}$$

$$\pi \sim p$$

Then:

$$\mathbb{E}_p[f^*(\pi^*) - f^*(\pi)]$$

$$\leq \frac{K}{\sigma} + \sigma \mathbb{E}[(\hat{f}(\pi) - f^*(\pi))^2]$$

PF Thm:

By Lemma

$$\mathbb{E} \text{Reg}_{CB} = \mathbb{E} \sum_{t=1}^T (f^*(x_t, \pi_t^*) - f^*(x_t, \pi_t))$$

$$\leq \sum_{t=1}^T \frac{K}{\gamma} + \gamma \sum_{t=1}^T \mathbb{E} (f^*(x_t, \pi_t^*) - f^*(x_t, \pi_t))^2$$

$$= \frac{KT}{\gamma} + \gamma \text{Reg}_{SQ} \quad \text{let } \gamma = \sqrt{\frac{KT}{\text{Reg}_{SQ}}} \leq 2\sqrt{KT \text{Reg}_{SQ}} \quad \square$$

Pf of Lemma:

$$\mathbb{E}_{\pi \sim P} [f^*(\pi^*) - f^*(\pi)]$$

(I)

$$= \mathbb{E}_P \left[(f^*(\pi^*) - \hat{f}(\pi^*)) \right. \\ \left. + (\hat{f}(\pi^*) - \hat{f}(\hat{\pi})) \right]$$

(II)

$$+ (\hat{f}(\hat{\pi}) - \hat{f}(\pi))$$

(III)

$$+ (\hat{f}(\pi) - f^*(\pi)) \left. \right]$$

(IV)

$\textcircled{\text{III}}$ \hookrightarrow $\textcircled{\text{IV}}$ easy to bound. Interesting part is $\textcircled{\text{I}}$ & $\textcircled{\text{II}}$.

$\textcircled{\text{I}}$

$$f^*(\pi^*) - \hat{f}(\pi^*)$$

$$= \sqrt{\frac{\partial p(\pi^*)}{\partial p(\pi^*)}} (f^*(\pi^*) - \hat{f}(\pi^*))$$

$$\leq \frac{\partial p(\pi^*)}{2} (f^*(\pi^*) - \hat{f}(\pi^*))^2$$

$$+ \frac{1}{2\partial p(\pi^*)} \leq \frac{\partial \mathbb{E}_p[(f^*(\pi) - \hat{f}(\pi))^2]}{2} + \frac{1}{2\partial p(\pi^*)}$$

① + ②

$$\leq \frac{\sigma}{2} \mathbb{E}_{\pi \sim p} \left[(f^*(\pi) - \hat{f}(\pi))^2 \right]$$

$$+ \frac{1}{2\sigma p(\pi^*)} + \left(\hat{f}(\pi^*) - \hat{f}\left(\frac{\pi}{2}\right) \right)$$

Recall: $\frac{1}{2\sigma p(\pi^*)} = \frac{\lambda + 2\sigma \left(\hat{f}\left(\frac{\pi}{2}\right) - \hat{f}(\pi^*) \right)}{2\sigma}$

$$\leq \frac{\sigma}{2} \mathbb{E} \left[(f^*(\pi) - \hat{f}(\pi))^2 \right] + \frac{\lambda}{2\sigma}$$