CS861: Theoretical Foundations of Machine Learning

Lecture 18 - 10/16/2023

University of Wi

Lecture 18:

ont'd), K-armed Bandit Lower Bound

Lecturer: Kirthe

Scribed by: Michael Harding and Congwei Yang

ected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the instructor.

In this lecture, we will first upper bound the regret for UCB, providing gap-dependent and worst-case bounds. We will then start our discussion on proving lower bounds for K-armed bandits.

Algorithm 1 The Upper Confidence Bound Algorithm

Require: time horizon Email: tutorcs@163.com

for $t = 1, \dots, K$ do

 $A_t \leftarrow t$

 $X_t \sim \nu_t$

end for

for $t = K+1, \ldots, T$ QQ: 749389476 $A_t \leftarrow \arg\max_{i \in [K]} \left(\widehat{\mu_{i,t-1}} + e_{i,t-1} \right)$

▶ Break ties arbitrarily

 $X_t \sim \nu_{A_t}$

end for

We will now present the theorem for the risk upper bounds for the UCB theorem once again, and pick up the proof where we left off.

Theorem 1 (UCB Risk Upper Bound). Let $\mathcal{P} = \{ \nu = \{ \nu_i \}_{i=1}^K : \nu_i \text{ } \sigma\text{-s}G, \ \mathbb{E}_{X \sim \nu_i}[X] \in [0,1] \ \forall \ i \in [K] \} \ be \ the$ class of σ -sub-Gaussian K-armed bandit models with means in [0,1]. Let $\mu_i := \mathbb{E}_{X \sim \nu_i}[X], \ \mu_* := \max_{i \in [K]} \mu_i,$ and denote $\Delta_i := \mu_* - \mu_i$. Then

$$R_T(\nu) \le 3K + \sum_{i:\Delta_i > 0} \frac{24\sigma^2 \log(T)}{\Delta_i} \tag{1}$$

$$\sup_{\nu \in \mathcal{P}} R_T(\nu) \le 3K + \sigma \sqrt{96KT \log(T)} \tag{2}$$

Proof As before, WLOG, we begin by letting $1 \ge \mu_1 \ge \cdots \ge \mu_K \ge 0$ for ease of notation. Also, we again define our good events

$$G_1 := \bigcap_{t > K} \left\{ \mu_1 < \hat{\mu}_{1,t} + e_{1,t} \right\}$$

$$G_1 := \bigcap_{t>K} \left\{ \mu_1 < \hat{\mu}_{1,t} + e_{1,t} \right\}$$
$$G_i := \bigcap_{t>K} \left\{ \mu_i > \hat{\mu}_{i,t} - e_{i,t} \right\}$$

At the end of our previous class, we proved that $\mathbb{P}(G_1^c)$, $\mathbb{P}(G_i^c) \leq \frac{1}{T}$ (we directly showed this for the case of G_1^c , remarking that the case for G_1^c is nearly identical). We will now show that $N_{i,t} := \sum_{s=1}^t \mathbb{I}_{\{A_s = i\}}$ is small for sub-optimal arms (a provided in the case for G_1^c is nearly identical). To show this, suppose arm i was last pulled on round t+1, where $t\geq 1$

$$(+e_{j,t}) \leftarrow \text{UCB Alg. construction}$$

and under G_i , we also \mathfrak{k}

$$\begin{array}{c} \mu_1 < \mu_i + 2e_{i,t} \Rightarrow \frac{\Delta_i}{2} < e_{i,t} = \sigma \sqrt{\frac{2\log(T^2t)}{N_{i,t}}} \\ \textbf{WeChat:} \underbrace{\textbf{CStut}}_{\Delta_i^2} \underbrace{\textbf{CT}}_{T>t} \\ \Rightarrow N_{i,t} < \underbrace{\frac{\Delta_i}{\Delta_i^2}}_{T>t} \leftarrow T>t \end{array}$$

Assign $\overrightarrow{m}_{e}^{N_{i,T}} = N_{i} P_{ro}^{1 \leq \frac{24\sigma^{2} \log(T)}{1}} E_{xam}^{1} Help$

Now, combining these results, we can write,

$$\mathbb{E}[N_{i,t}] = \underbrace{\mathbb{E}[N_{i,t}|G_1 \cap G_i] \mathbb{P}(G_1 \cap G_i)}_{\leq \underbrace{\mathsf{TUTOTCS}}} + \mathbb{E}[N_{i,t}|G_1^c \cup G_i^c] \mathbb{P}(G_1^c \cup G_i^c) \leq 3 + \underbrace{\frac{24\sigma^2 \log(T)}{\Delta_i^2}}_{\Delta_i^2}$$

Then, by the regret decomposition result shown towards the end of last class, we can write,

$$Q_{R,Q}:=\sum_{i:\Delta_{i}>0}49389476_{i:\Delta_{i}>0}\frac{24\sigma^{2}\log(T)}{\Delta_{i}},$$

where we leverage the fact that $\Delta_i \in [0,1]$ and there are at most K-1 summands. This proves the gap-dependent bound in (1) Tritle springer broken can choose some value $\Delta > 0$ and rewrite our result above as thus:

$$R_{T}(\nu) = \sum_{i:\Delta_{i}>0} \Delta_{i} \mathbb{E}[N_{i,t}]$$

$$= \sum_{i:\Delta_{i}\in(0,\Delta]} \Delta_{i} \mathbb{E}[N_{i,t}] + \sum_{i:\Delta_{i}>\Delta} \Delta_{i} \mathbb{E}[N_{i,t}]$$

$$\leq \Delta \sum_{i:\Delta_{i}\in(0,\Delta]} \mathbb{E}[N_{i,t}] + \sum_{i:\Delta_{i}>\Delta} \frac{24\sigma^{2}\log(T)}{\Delta} + 3K$$

$$\leq 3K + \Delta T + \frac{24\sigma^{2}\log(T)}{\Delta}$$

Then, because this holds for all $\Delta > 0$, we are free to optimize over values of Δ , giving us in particular $\Delta = \sigma \sqrt{\frac{24K \log(T)}{T}}$. Therefore,

$$R_T(\nu) \le 3K + \sigma \sqrt{96KT\log(T)}$$

and because this result holds for all $\nu \in \mathcal{P}$, and the bound has no dependence on ν , then we can write,

$$\sup_{\nu \in \mathcal{P}} R_T(\nu) \le 3K + \sigma \sqrt{96KT \log(T)},$$

which is exactly the statement in (2).

the gap-independent bound. We will use similar techniques Next, we will presen for linear bandits in sub

Alternative 🖥

We will first decompose

¶



$$R_T = \mathbb{E}\left[\sum_{t=1}^T (\mu_* - X_t)\right]$$

WeChat: Estutoros

$$R_T = \sum_{t=1}^{\infty} a_{\mu 1} \sum_{\mu_{A_t}} u_t \text{ or } c_{\mu} \text{ or } c_{\mu_1} \text{ o$$

Note we have P(G) $T_{t=1}^T$ T_t $T_$

Claim: Under the event G, $\mu_1 - \mu_{A_t} \leq 2e_{A_{t},t-1}$

- If A_t is an optimal arm, then $\mu_1 \mu_{A_t} \leq 0 \leq 2e_{A_t,t-1}$. If not, $\mu_1 \leq \hat{\mu}_{1,t-1}$ is $\hat{\mu}_{1,t-1}$ in $\hat{\mu}_{1,t-1}$ in $\hat{\mu}_{1,t-1}$ in $\hat{\mu}_{1,t-1}$ is under $\hat{\mu}_{1,t-1}$ where the first inequality is under $\hat{\mu}_{1,t-1}$ and the last inequality is under $\hat{\mu}_{1:\Delta_t>0}$ $\hat{\mu}_{1:\Delta_t>0}$ $\hat{\mu}_{1:\Delta_t>0}$ $\hat{\mu}_{1:\Delta_t>0}$

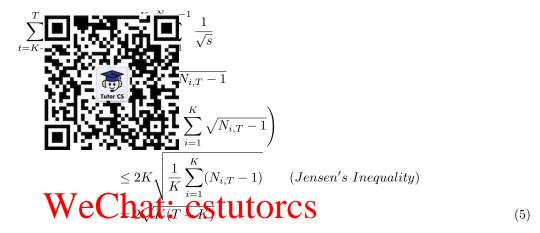
Then,

$$\sum_{t=1}^{T} (\mu_1 - \mu_{A_t}) \le K + \sum_{t=K+1}^{T} 2\sigma \sqrt{\frac{2\log(1/\delta_t)}{N_{A_t, t-1}}}$$

$$\le K + \sum_{t=K+1}^{T} 2\sigma \sqrt{\frac{2\log(T^2 t)}{N_{A_t, t-1}}}$$

$$\le K + \sigma \sqrt{24\log(T)} \sum_{t=K+1}^{T} \frac{1}{\sqrt{N_{A_t, t-1}}}$$
(4)

We will now focus on the last summation:



Here the first inequality follows from $\sum_{s=1}^{m} \frac{1}{\sqrt{s}} \leq 2\sqrt{m}$, which we have proved below.

Combining (3), (4), (5), we obtain $R_T < 2K + \sigma \sqrt{96KT \log T}$ To prove, $\sum_{s=1}^{m} \frac{1}{\sqrt{s}} \le 2\sqrt{m}$, we will bound the sum of a decleasing function by an integral as follows: $\sum_{s=1}^{m} \frac{1}{\sqrt{s}} \le \int_0^m \frac{1}{\sqrt{s}} ds = (2s^{1/2})|_0^m = 2\sqrt{m}$.

2 K-armed bands alwertutores @ 163.com

In this section, we will prove the following lower bound on the minimax regret: $\inf_{\Pi} \sup_{\nu \in \mathcal{P}} R_T(\Pi, \nu) \in \Omega(\sqrt{KT})$. To do so, recall the following results 3 each proof of Le Cam's method (Lecture 9, Lemma 1 and Corollary 1).

Lemma 1. Let P_0 , P_1 be two distributions and A be any event. Then,

$$\begin{array}{c|c} \mathbf{httpS.} / 2 & P_1(A^c) / 2 & P_2(A^c) \\ & & \mathbf{httpS.} / 2 & \mathbf{httpS.} / 2 & \mathbf{httpS.} \\ & \geq \frac{1}{2} \exp(-KL(P_0, P_1)) \end{array}$$

When applying this inequality, the KL divergence will be between distributions of action-reward sequences $A_1, X_1, \dots, A_T, X_T$ induced by the interaction of a policy π with different bandit models. The following lemma will be helpful in computing the KL divergence.

Lemma 2 (KL divergence decomposition). Let ν , ν' be two K-armed bandits models. For a fixed policy Π , let P, P' denote the probability distribution over the sequence of actions and rewards $A_1, X_1, \dots, A_T, X_T$ under ν , ν' , respectively. Let \mathbb{E}_{ν} denote the expectation under bandit model ν . Then $\forall T \geq 1$,

$$KL(P, P') = \sum_{i=1}^{K} \mathbb{E}_{\nu}[N_{i,T}]KL(\nu_i, \nu_i')$$

where $N_{i,T} = \sum_{t=1}^{T} \mathbf{1}_{\{A_t = i\}}$

Intuitively, suppose we pulled arm 1 N_1 times. As the observations are independent $KL(P,P') = N_1KL(\nu_1,\nu_1')$. Next, consider a nonadaptive policy which pulls arm i N_i times for $i=1,\cdots,K$. We then have $KL(P,P') = \sum_{i=1}^{K} N_iKL(\nu_i,\nu_i')$. The above lemma says that a similar result holds when we use an adaptive policy.

Proof Proof of Lemma 2 Consider any given sequence $a_1, x_1, \dots, a_T, x_T$. Let p, p' denote the Radon-Nikodym derivatives of $\underline{P, P'}$ respectively. Let $\tilde{\nu}_i, \tilde{\nu}'_i$ denote the Radon-Nikodym derivatives of ν_i, ν'_i , respectively.

Consider for fixed a

$$\prod_{t \text{ Totor CS}} p(a_t, x_t \mid a_1, x_1, \cdots, a_{t-1}, x_{t-1})$$

$$\underset{\log}{Assignment} \underbrace{Project}_{p'(a_1,x_1,\cdots,a_t,x_t)} \underbrace{Project}_{\tilde{\nu}_{a_1}(x_1)\cdots\tilde{\nu}_{a_t}(x_T)} \underbrace{Fxam}_{Help}$$

Email: tutores @ B.com

To be continued next lecture...

QQ: 749389476

https://tutorcs.com

5