

Lecture 19:  $K$ -armed bandit lower bounds, generalized linear bandits

Lecturer: Kirthevasan Kandasamy

Scribed by: Alex Clinton, Yamin Zhou

**Disclaimer:** These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the instructor.

In the previous lecture we proved an upper bound for the UCB and began analyzing a lower bound. In this lecture, we will continue the proof of that lower bound which is gap independent. We will then also provide a gap dependent lower bound for  $k$ -armed bandits. Finally, we will introduce a structured bandit model which generalizes the  $K$ -armed setting.

## 1 $K$ -armed bandits lower bounds

**Lemma 1.** Let  $\nu, \nu'$  be two bandit models, and  $P, P'$  bet the prob distribution of  $A_1, x_1, \dots, A_i, x_i$  due to the interaction of a policy  $\pi$  with  $\nu, \nu'$  respectively, then

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

**Proof** we showed from the previous lecture, for any sequence  $a_1, x_1, \dots, a_T, x_T$

$$\log \left( \frac{p(a_1, x_1, \dots, a_T, x_T)}{p'(a_1, x_1, \dots, a_T, x_T)} \right) = \sum_{t=1}^T \log \left( \frac{\hat{\nu}_{A_t}(x_t)}{\hat{\nu}'_{A_t}(x_t)} \right)$$

Therefore, we have

$$\begin{aligned} KL(P, P') &= \mathbb{E}_p \left[ \log \left( \frac{P(A_1, x_1, \dots, A_i, x_i)}{P'(A_1, x_1, \dots, A_T, x_T)} \right) \right] \\ &= \mathbb{E}_p \left[ \sum_{t=1}^T \log \left( \frac{\hat{\nu}_{A_t}(x_t)}{\hat{\nu}'_{A_t}(x_t)} \right) \right] \\ &= \sum_{t=1}^T \mathbb{E}_p \left[ \log \left( \frac{\nu_{A_t}(x_t)}{\nu'_{A_t}(x_t)} \right) \sum_{i=1}^K \mathbb{1}(A_t = i) \right] \\ &= \sum_{i=1}^K \sum_{t=1}^T \mathbb{E}_p \left[ \underbrace{\mathbb{E}_p \left[ \log \left( \frac{\nu_{A_t}(x_t)}{\nu'_{A_t}(x_t)} \right) \mathbb{1}(A_t = i) \mid A_t \right]}_* \right] \\ &= \sum_{i=1}^K KL(\nu_i, \nu'_i) \left( \sum_{t=1}^T \mathbb{E}[\mathbb{1}(A_t = i)] \right) \quad \text{from the result of } * \text{ and } KL \text{ is not related to } t \\ &= \sum_{i=1}^K KL(\nu_i, \nu'_i) \mathbb{E}[N_{i,T}] \end{aligned}$$

where

$$\begin{aligned}
* &= \mathbb{1}(A_t = i) \mathbb{E}_p \left[ \log \left( \frac{\nu_{A_t}(x_t)}{\nu'_{A_t}(x_t)} \right) \middle| A_t \right] \\
&= \mathbb{1}(A_t = i) KL(\nu_{A_t}, \nu'_{A_t}) \\
&= \mathbb{1}(A_t = i) KL(\nu_i, \nu'_i)
\end{aligned}$$

□

This KL divergence decomposition lemma now allows us to state the following minimax lower bound for  $k$ -armed bandits. The idea of the proof is that given any policy  $\pi$ , we can construct two bandit instances such that  $\pi$  will perform poorly in one of them.

**Theorem 1.** (*Minimax lower bound for  $k$ -armed bandits*) Let  $P = \{\nu = \nu_i, i \in [K]\}$ ,  $\nu_i$  is  $\sigma$  subGaussian for all  $i \in [K]$ . Then, if  $K > 1$ , for some universal constant  $C$ ,

$$\inf_{\pi} \sup_{\nu \in P} R_T(\pi, \nu) \geq C\sigma\sqrt{T(K-1)}$$

**Proof** Let  $\pi$  be given, consider the following two bandit models  $\nu, \nu'$  (constructed based on  $\pi$ ) as follows:

- Let  $\nu = \{\nu_i = N(\mu_i, \sigma^2)_{i \in [K]}\}$ , where  $\mu_1 = \delta$ , and  $\mu_i = 0, \forall i \in 2, \dots, K$ ,  $\mu = (\delta, 0, 0, \dots, 0)$ . Here,  $\delta > 0$  is a value we will specify shortly.
- Let  $\mathbb{E}_\nu$  denote the expectation with respect to the sequence  $A_1, x_1, \dots, A_T, x_T$  due to  $\pi$ 's interaction with  $\nu$ . Since  $\sum_{i=1}^k \mathbb{E}_\nu[N_{i,t}] = T$ ,  $\exists$  some  $j \in 2, \dots, k$  s.t  $\mathbb{E}_\nu[N_{j,t}] \leq \frac{T}{k-1}$
- Let  $\nu' = \{\nu_i = N(\mu'_i, \sigma^2)\}_{i \in [k]}$  where

$$\mu'_i = \begin{cases} \mu_i & \text{if } i \neq j \\ 2\delta & \text{if } i = j \end{cases}$$

Therefore,  $\mu' = (\delta, 0, 0, \dots, \underbrace{2\delta}_j, 0, \dots, 0)$

- Let  $P, P'$  denote the prob distributions of  $A_1, x_1, \dots, A_T, x_T$  due to  $\pi$ 's interaction with  $\nu, \nu'$ .

We know,

$$\begin{aligned}
R_T(\pi, \nu) &\geq P \left( N_{1,T} \leq \frac{T}{2} \right) \frac{T\delta}{2} \\
R_T(\pi, \nu') &\geq P' \left( N_{j,T} \leq \frac{T}{2} \right) \frac{T\delta}{2} \geq P' \left( N_{1,T} \geq \frac{T}{2} \right) \frac{T\delta}{2}
\end{aligned}$$

Now, note that we can write

$$\sup_{\nu \in P} R_T(\pi, \nu) \geq \max(R_T(\nu, \pi), R_T(\nu', \pi)) \geq \frac{1}{2} \underbrace{(R_T(\nu, \pi) + R_T(\nu', \pi))}_{(*)}$$

We will now bound the term  $(\star)$  as follows

$$\begin{aligned}
(\star) &\geq \frac{T\delta}{2} \left( P(N_{1,T} \leq \frac{T}{2}) + P'(N_{1,T} > \frac{T}{2}) \right) && \text{from the definition} \\
&\geq \frac{T\delta}{4} \exp(-KL(p, p')) && \text{from the Le Cam's lemma} \\
&= \frac{T\delta}{4} \exp\left(\mathbb{E}_\nu[N_{j,T}] \frac{(2\delta)^2}{2\sigma^2}\right) && \text{Divergece Decomposition + KL of Gaussian} \\
&\geq \frac{T\delta}{4} \exp\left(-\frac{T}{K-1} \frac{2\delta^2}{\sigma^2}\right)
\end{aligned}$$

Now, if we choose  $\delta = \sigma \sqrt{\frac{K-1}{T}}$ , Then, we are able to get  $\star \geq (\frac{1}{4}e^{-2})\sqrt{T(K-1)}$ , and hence  $\inf_\pi \sup_{\nu \in \mathcal{P}} R_T(\pi, \nu) \geq C\sigma\sqrt{T(K-1)}$ .  $\square$

## 2 Gap-independent lower bounds

Recall from our analysis of UCB that there are two types of bounds: gap dependent bounds (those that depend on  $\Delta_i$ ) and gap independent bounds (those that do not depend on  $\Delta_i$ ). In addition to the gap independent lower bound we just proved, there is also a gap dependent lower bound for  $k$ -armed bandits. Although we will not prove it here, this lower bound is given by the following theorem.

**Theorem 2.** (Theorem 16.4 in LS)

Let  $\nu$  be a given  $K$ -armed bandit model with  $\sigma$ -sub Gaussian rewards. Let  $\mu = \mu(\nu)$  be the means of the arms. Let  $\mathcal{P}(\nu) = \{\nu' : \mu_i(\nu) \in [\mu_i, \mu_i + 2\Delta_i], v_i \text{ is } \sigma\text{-sub Gaussian}\}$ . Say  $\pi$  is a policy such that  $R_T(\pi, \nu') \leq cT^p, \forall \nu' \in \mathcal{P}(\nu)$  for some  $c > 0$  and  $p \in (0, 1)$ . Then

$$R_T(\pi, \nu) \geq \frac{1}{2} \sum_{i:\Delta_i > 0} \frac{\left((1-p)\log(T) + \log\left(\frac{\Delta_i}{8c}\right)\right) \sigma^2}{\Delta_i}$$

**Remark** At a high level this theorem says that if a policy does well on all “similar” problems then it does at least as poorly as the given expression on the original problem.

## 3 Stochastic bandits in a generalized linear model

One potential criticism of the bandit model we have studied thus far is its restriction of the action space to  $K$  specific choices. In the generalized linear bandit model we are about to introduce, we will allow for an infinite action space, but assume additional structure on the rewards.

**Definition 1.** A generalized linear bandit model consists of the following components:

1. Action space  $\mathcal{A} \subseteq [-1, 1]^d$ . For reasons of convenience which will become clear shortly, we will assume that the basis vectors  $e_1, \dots, e_d$  are in  $\mathcal{A}$ .
2. Parameter space  $\Theta \subseteq [-1, 1]^d$
3. True parameter (unknown)  $\theta_* \in \Theta$
4. When we choose an action (arm)  $A_t$ , we observe  $X_t = f(\theta_*^T A_t) + \varepsilon_t$  where  $\mathbb{E}[\varepsilon_t] = 0$  and  $\varepsilon_t$  is  $\sigma$ -sub Gaussian. Here  $\varepsilon_t$  can be thought of as noise.

5. Here  $f$  is known and has the following properties

- (a)  $f$  is strictly increasing with  $f'(x) \geq c > 0$
- (b)  $f$  is  $L$ -Lipschitz
- (c)  $f'$  is continuous

**Example** ( $f$  is the identity in a linear bandit model): In this case we have the following expression for the regret,  $R_T = Tf(\theta_*^T a_*) - \mathbb{E} \left[ \sum_{t=1}^T X_t \right]$ , where  $a_* = \arg \max_{a \in \mathcal{A}} f(\theta_*^T a)$ . Notice that this allows us to recover our original bandit model and so our new framework is broader than how we first introduced  $k$ -armed bandits.

### 3.1 A UCB algorithm (Based on Filippi et al. 2010)

In order to get a UCB algorithm, we need to know how to estimate  $\theta_*$  from data and construct UCBs. We define the following quantities:

- $\hat{\theta}_t = \arg \min_{\theta \in \Theta} \left\| \sum_{s=1}^t A_s (f(A_s^T \theta) - X_s) \right\|_{V_t^{-1}}$  where  $v_t = \sum_{s=1}^t A_s A_s^T$  and  $\|y\|_Q^2 = y^T Q y$  (here  $Q$  must be positive semi-definite)
- $\rho(t) = \frac{2L\sigma}{c} \sqrt{(3 + 2 \log(1 + 2d)) \cdot 2d \log(t) \log(dT^2)} \in \tilde{O}(\sqrt{d})$

In the following algorithm, for the first  $d$  rounds, we pull each basis vector which is analogous to pulling each of the  $k$  arms once at the start of our original UCB algorithm. For the remaining rounds we then pull the arm with the highest new confidence bound which is again analogous the original UCB algorithm where we pull the arm with the highest original confidence bound.

---

#### Algorithm 1 UCB

---

**Require:** a time horizon  $T$

**for**  $t = 1, \dots, d$  **do**

Choose  $A_t = e_t$  (the  $t^{\text{th}}$  basis vector)

**end for**

**for**  $t = d + 1, \dots, T$  **do**

Choose  $A_t = \arg \max_{a \in \mathcal{A}} \underbrace{f(\hat{\theta}_t^T a)}_{\text{exploitation}} + \underbrace{\rho(t) \|a\|_{V_{t-1}^{-1}}}_{\text{exploration}}$

**end for**

---