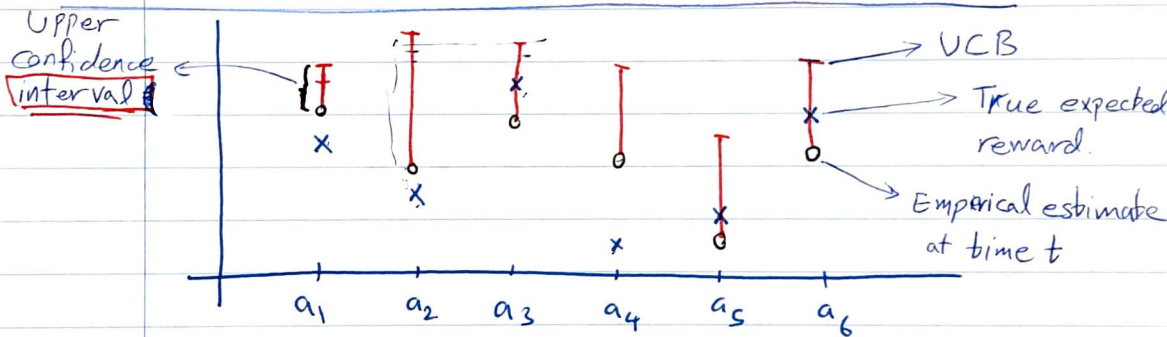


## CS 59000 RL

## Finite-armed bandit

Agenda:

- UCB
- Lower bound.
- Minimax optimal
- Some more side information
- Adversarial bandit.



Recall the UCB type algorithm we talk in the previous lecture.

Theorem: UCB algorithm on a stochastic  $K$ -armed bandit with 1-sub-Gaussian reward, achieves regret of

$$R_n(\pi, \nu) \leq 3 \sum_{i=1}^K \Delta_i + \sum_{i: \Delta_i > 0} \frac{16 \log(n)}{\Delta_i}$$

$\downarrow$  Policy       $\downarrow$  environment  
 UCB algorithm

when  $\delta = \frac{1}{n^2}$

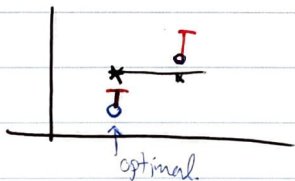
Proof: Recall that  $R_n = \sum_{i=1}^K \Delta_i E[T_i(n)]$   
we construct proof

by bounding  $E[T_i(n)]$  for each sub optimal arm

When we are sure that we are not mistakenly  
pulling arm  $i > 1$  ? (For simplicity, assume  
arm 1 is one of the  
optimal arms)

Consider the following events.

$\mu_1 < \min_{t \in [n]} \text{UCB}_1(t, \delta)$ ; Upper Confidence  
bound on the optimal  
arm is a valid  
bound for  $1 \leq t \leq n$



• For some number  $u_i$ :

$$\hat{\mu}_{i, u_i} + \sqrt{\frac{2}{u_i} \log \frac{1}{\delta}} < \mu_1$$

empirical mean after pulling  
i'th arm for  $u_i$  times.

After  $u_i$  times  
of pulling i'th arm  
the upper confidence  
bound is smaller  $\mu_1$ .

Consider the event  $G_i$ :

$$G_i = \left\{ \mu_1 < \min_{t \in [n]} \text{UCB}_1(t, \delta) \cap \left( \hat{\mu}_{i, u_i} + \sqrt{\frac{2}{u_i} \log \frac{1}{\delta}} < \mu_1 \right) \right\}$$

First we show that in event  $G_i$ ,  $T_i(n) < u_i$  using contradiction.

If  $T_i(n) > u_i$ , then there is a  $t \leq n$  where  $T_i(t-1) = u_i$  and  $A_t = i$ .

By definition of  $UCB_i(t-1, \delta)$ , we have

$$\begin{aligned} UCB_i(t-1, \delta) &= \hat{\mu}_i(t-1) + \sqrt{\frac{2 \log(\frac{1}{\delta})}{T_i(t-1)}} \\ &= \hat{\mu}_{i, u_i} + \sqrt{\frac{2 \log(\frac{1}{\delta})}{u_i}} \end{aligned}$$

$$\begin{array}{l} \text{under event } G_i \\ \xrightarrow{\quad} < \mu_i \\ \searrow < UCB_i(t-1, \delta) \end{array}$$

therefore at time  $t$ , we would not choose  $i$ th arm  $\rightarrow$  contradiction  $\Rightarrow$  Therefore  $T_i(n) \leq u_i$ .

— Using this inequality we have:

$$\begin{aligned} E[T_i(n)] &= \int_{G_i} T_i(n) d\mathbb{P}_{\mathcal{V}, \pi} + \int_{G_i^c} T_i(n) d\mathbb{P}_{\mathcal{V}, \pi} \\ &= E[I_{G_i} T_i(n)] + E[I_{G_i^c} T_i(n)] \end{aligned}$$

$$\Rightarrow E[T_i(n)] \leq E[I_{G_i} u_i] + n E[I_{G_i^c}]$$

$$= u_i \underbrace{E[I_{G_i}]}_{P(G_i)} + n P(G_i^c)$$

$$\leq \underbrace{u_i}_{\downarrow} + n \underbrace{P(G_i^c)}$$

we will choose  $u_i$   
later.

Let us upper bound this.

$$G_i^c = \left\{ \mu_i \geq \min_{t \in [n]} \text{UCB}_1(t, \delta) \right\}$$

$$\vee \left\{ \mu_i, u_i + \sqrt{\frac{2 \log \frac{1}{\delta}}{u_i}} \geq \mu_i \right\}$$

Let's study the first event

$$\left\{ \mu_i \geq \min_{t \in [n]} \text{UCB}_1(t, \delta) \right\}$$

$$\subset \left\{ \mu_i \geq \min_{s \in [n]} \hat{\mu}_{1,s} + \sqrt{\frac{2 \log \frac{1}{\delta}}{s}} \right\}$$

$$\subset \bigcup_{s \in [n]} \left\{ \mu_i \geq \hat{\mu}_{1,s} + \sqrt{\frac{2 \log \frac{1}{\delta}}{s}} \right\}$$

$\Rightarrow$  using union bound  $\rightarrow$

$$P\left(\mu_i \geq \min_{t \in [n]} \text{UCB}_1(t)\right) \leq \sum_{s=1}^n P\left(\mu_i \geq \hat{\mu}_{1,s} + \sqrt{\frac{2 \log \frac{1}{\delta}}{s}}\right) \leq n\delta$$

Page 5

Now let's study  $\left\{ \hat{\mu}_i, u_i + \sqrt{\frac{2 \log \frac{1}{\delta}}{u_i}} \geq \mu_i \right\}$

Let  $u_i$  be large enough such that

$$\Delta_i - \sqrt{\frac{2 \log \frac{1}{\delta}}{u_i}} \geq C \Delta_i$$

$$\begin{aligned} \Rightarrow \mathbb{P}\left(\hat{\mu}_i, u_i + \sqrt{\frac{2 \log \frac{1}{\delta}}{u_i}} \geq \mu_i\right) &= \mathbb{P}\left(\hat{\mu}_i, u_i - \mu_i \geq \Delta_i - \sqrt{\frac{2 \log \frac{1}{\delta}}{u_i}}\right) \\ &\leq \mathbb{P}\left(\mu_i, u_i - \mu_i \geq C \Delta_i\right) \leq \exp\left(-\frac{u_i C^2 \Delta_i^2}{2}\right) \end{aligned}$$

using union bound

$$\Rightarrow \mathbb{P}(G_i^c) = n\delta + \exp\left(-\frac{u_i C^2 \Delta_i^2}{2}\right)$$

$\Rightarrow$  therefore

$$\mathbb{E}[T_i(n)] \leq u_i + n \left( n\delta + \exp\left(-\frac{u_i C^2 \Delta_i^2}{2}\right) \right)$$

Let us set  $u_i = \left\lceil \frac{2 \log(\frac{1}{\delta})}{(1-c)^2 \Delta_i^2} \right\rceil$  and set  $\delta = \frac{1}{n^2}$

$$\mathbb{E}[T_i(n)] \leq \frac{2 \log \frac{1}{\delta}}{(1-c)^2 \Delta_i^2} + 1 + n \frac{1-2c^2}{1-c}$$

Let's set  $C = \frac{1}{2} \Rightarrow$  then

$$\mathbb{E}[T_i(n)] \leq 3 + \frac{16 \log(n)}{\Delta_i^2}$$



$$R_n = \sum \Delta_i E[T_i(n)]$$

$$\Rightarrow R_n \leq 3 \sum_{i=1}^K \Delta_i + \sum_{\substack{i=1 \\ \Delta_i > 0}}^K \frac{16 \log(n)}{\Delta_i}$$

Theorem: The regret of UCB algorithm on  $L$ -sub-Gaussian stochastic  $K$ -armed bandit is upper bounded by

$$R_n \leq 8 \sqrt{n K \log(n)} + 3 \sum_{i=1}^K \Delta_i$$

proof: For a parameter  $\Delta > 0$

$$R_n = \sum_{i=1}^K \Delta_i E[T_i(n)]$$

$$= \sum_{i: \Delta_i < \Delta} \Delta_i E[T_i(n)] + \sum_{i: \Delta_i \geq \Delta} \Delta_i E[T_i(n)]$$

$$\leq n\Delta + \frac{16 K \log(n)}{\Delta} + 3 \sum_i \Delta_i$$

$$\Rightarrow \text{Let's set } \Delta = \sqrt{\frac{16 K \log n}{n}}$$

$$\Rightarrow R_n \leq 8 \sqrt{n K \log n} + 3 \sum_{i=1}^K \Delta_i$$

$$\frac{R_n}{\log n} \rightarrow$$

Lower Bound: The worst case regret of a policy  $\pi$  on a class of environment  $\mathcal{E}$  is:

$$R_n(\pi, \mathcal{E}) = \sup_{v \in \mathcal{E}} R_n(\pi, v)$$

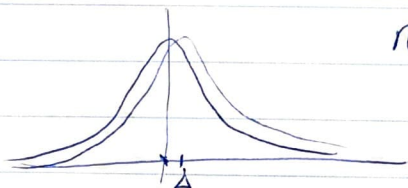
Let  $\Pi$  be the set of all policies, The **minimax** regret of  $\Pi$  on  $\mathcal{E}$  is

$$R_n^*(\Pi, \mathcal{E}) = \inf_{\pi \in \Pi} R_n(\pi, \mathcal{E}) = \inf_{\pi \in \Pi} \sup_{v \in \mathcal{E}} R_n(\pi, v)$$

A policy is minimax optimal if it achieves minimax regret

Theorem. Let  $\mathcal{E}^K$  be the set of  $K$ -armed stochastic Gaussian bandits, with unit variance and mean  $\mu \in [0, 1]^K$ . Then there exist a constant  $c > 0$  such that for  $K > 1$  and  $n \geq K$ , it holds

$$R_n^*(\mathcal{E}^K) \geq c \sqrt{(K-1)n}$$



$n$

$\mu_1(\Delta, \dots, \Delta, \dots)$

$\mu(\Delta, \dots, 2\Delta, \dots)$