

## LECTURE 20

# THEORETICAL ANALYSIS OF BANDIT ALGORITHMS

Setup:  $K$  coins

success probabilities:  $p_1, \dots, p_K$

$T$  rounds. (e.g.:  $K=9, p_1=0.1, p_2=0.2, \dots, p_K=0.9$ )

Each round: You pick a coin for tossing  
1 \$ if heads  
\$0 if tails

Given an "algorithm";

Reward ( $t$ ) = Money made after  $t$  rounds.

↑  
number of rewards

$$\text{Regret}(t) = 0.9t - \text{Reward}(t)$$
$$= t \cdot \max_a p_a - \text{Reward}(t)$$

Averaged Regret ( $t$ )

$$= \mathbb{E}(\text{Regret}(t))$$
$$= t \cdot \max_a p_a - \mathbb{E}(\text{Reward}(t))$$

Averaged Regret increases with  $t$

Averaged Regret ( $t+1$ )

$$\begin{aligned}
&= (t+1) \max_a p_a - \mathbb{E}[\text{Reward}(t+1)] \\
&= t \max_a p_a + \mathbb{E}(\text{Reward}(t)) \\
&\quad + \underbrace{\max_a p_a - \mathbb{E}[\text{Reward}(t+1) - \text{Reward}(t)]}_{\geq 0} \\
&\quad \underbrace{\max_a p_a - \mathbb{E}[\text{payoff in Round}(t+1)]}_{\geq 0}
\end{aligned}$$

So Averaged Regret  $(t+1) \geq$  Averaged Regret  $(t)$

---

Example ①  $K = 9$

$$p_1 = 0.1, p_2 = 0.2, \dots, p_8 = 0.8, p_9 = 0.9$$

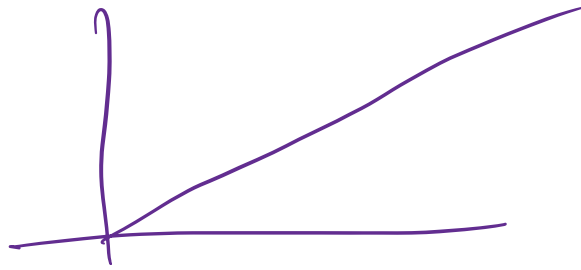
Algorithm: At each round, picks one of coin 8 or coin 9 at random.

Qn: What is the Averaged Regret of this Algorithm?

Averaged Regret  $(t)$

$$\begin{aligned}
&= t \max_a p_a - \mathbb{E}(\text{Reward}_t) \\
&= t(0.9) - \mathbb{E} \sum_{s=1}^t \mathbb{I} \{ \text{success in } s^{\text{th}} \text{ round} \} \\
&= t(0.9) - \sum_{s=1}^t \mathbb{P} \{ \text{success in } s^{\text{th}} \text{ round} \} \\
&= t(0.9) - \sum_{s=1}^t \left[ \frac{1}{2}(0.8) + \frac{1}{2}(0.9) \right]
\end{aligned}$$

$$= t(0.9) - t(0.85) = t(0.9 - 0.85) \\ = t(0.05)$$



Linear Regret Algorithm

Example 2:

In round  $t$ :

0.8 coin  $\rightarrow p_t$

0.9 coin  $\rightarrow 1 - p_t$

$p_t$  : decreases with  $t$ .

(e.g.  $p_t = \frac{1}{t}$ )

Averaged Regret for this algorithm.

$$= t(0.9) - \mathbb{E} \sum_{s=1}^t \mathbb{I} \{ \text{success in Round } s \}$$

$$= t(0.9) - \sum_{s=1}^t \mathbb{P}(\text{success in Round } s)$$

$$= t(0.9) - \sum_{s=1}^t [(0.8)(p_s) + (0.9)(1 - p_s)]$$

$$= t(0.9) - \sum_{s=1}^t [0.9 - 0.1 p_s]$$

$$= 0.1 \sum_{s=1}^t p_s \quad \text{with } p_s = \frac{1}{s}$$

$$\text{Averaged Regret} \\ = 0.1 \left( 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{t} \right)$$

$$\sim (0.1) (\log t)$$

$$\text{If } p_s = \frac{1}{s^2}$$

$$\text{Averaged Regret} = 0.1 \left( \sum_{s=1}^t \frac{1}{s^2} \right) \leq \text{Constant.}$$

---

### Uniform Confidence Bound.

① Basic Exploration ( $m=1$ )

② In Round ( $t$ ),

$T_a(t)$  = # times coin  $a$  is picked in the first  $t$  rounds

$X_a(t)$  = # successes from coin  $a$  in the first  $t$  rounds.

$$UCB_a(t, \epsilon) = \frac{X_a(t)}{T_a(t)} + \sqrt{\frac{\log(1/\epsilon)}{2 T_a(t)}}$$

$$\swarrow UCB_a(t-1, \epsilon)$$

pick the coin  $a$  which maximizes this.

---

$$\text{Idea } \rightarrow \mathbb{P} \left[ p_a \leq \underbrace{UCB_a(t, \epsilon)} \right] \geq 1 - \epsilon$$

---

Hoeffding for Binomial

$$X \sim \text{Bin}(n, p), \quad t \geq np$$

$$\mathbb{P}(X \geq t) \leq \exp\left[-\frac{2(t-np)^2}{n}\right]$$

=  $\delta$

Solve for  $t$ :

$$t = np + \sqrt{\frac{n}{2} \log \frac{1}{\delta}}$$

$$\frac{t}{n} = p + \sqrt{\frac{\log 1/\delta}{2n}}$$

$$\mathbb{P}\left[X \geq np + \sqrt{\frac{n}{2} \log \frac{1}{\delta}}\right] \leq \delta$$

$$\mathbb{P}\left[p \leq \frac{X}{n} - \sqrt{\frac{\log 1/\delta}{2n}}\right] \leq \delta$$

$$\mathbb{P}\left[p > \frac{X}{n} - \sqrt{\frac{\log 1/\delta}{2n}}\right] \geq 1 - \delta$$

LCB (Lower Confidence Bound)

UCB from Hoeffding:

$$\mathbb{P}\left[p < \frac{X}{n} + \sqrt{\frac{\log 1/\delta}{2n}}\right] \geq 1 - \delta$$

UCB Bandit Algorithm

After  $t-1$  rounds, consider the maximum possible value for each  $p_a$  subject to a given prob  $\delta$ .

$$\boxed{UCB_a(t-1, s)}$$

$\mathbb{P}_a$

$$\rightarrow \begin{cases} T_a(t-1) \\ X_a(t-1) \end{cases}$$

$$\frac{X_a(t-1)}{T_a(t-1)} + \sqrt{\frac{\log 1/\delta}{2T_a(t-1)}}$$

Qn: What is the Averaged Regret of UCB Bandit Algorithm?

Ans: Logarithmic. ( $\delta_t := \frac{1}{t^3}$ )

Sketch of proof:

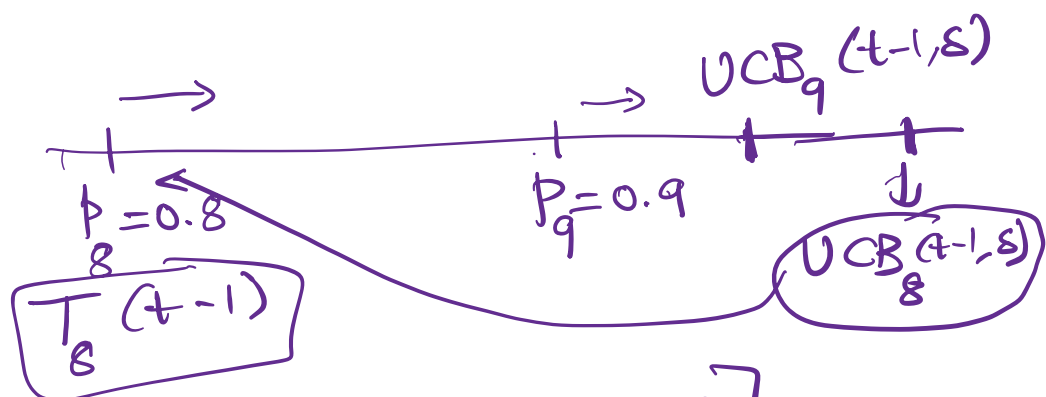
$$p_1 = 0.1, p_2 = 0.2, \dots, p_8 = 0.8, p_9 = 0.9$$

In round  $t$ , what are the chances of picking coin 8 over coin 9.

Intuition: This should be small when  $t$  is large.

$$\mathbb{P} \left[ UCB_8(t-1, \delta) > UCB_9(t-1, \delta) \right]$$

$$\mathbb{P} \left[ \frac{X_8(t-1)}{T_8(t-1)} + \sqrt{\frac{\log 1/\delta}{2T_8(t-1)}} > \frac{X_9(t-1)}{T_9(t-1)} + \sqrt{\frac{\log 1/\delta}{2T_9(t-1)}} \right]$$



$$\mathbb{P} \left[ \underbrace{UCB_8(t-1, \delta)}_{\geq 0.9} \right]$$

$$= \mathbb{P} \left[ \frac{X_8(t-1)}{\underbrace{T_8(t-1)}} + \sqrt{\frac{\log 1/\delta}{2 T_8(t-1)}} \geq 0.9 \right]$$

$$= \mathbb{P} \left[ \frac{X_8(t-1)}{T_8(t-1)} - 0.8 \geq \underbrace{0.9 - 0.8}_{\frac{\log 1/\delta}{\sqrt{2 T_8(t-1)}}} \right]$$

$$= \mathbb{P} \left[ \frac{X_8(t-1)}{T_8(t-1)} - 0.8 \geq \sqrt{\frac{\log 1/\delta}{2 T_8(t-1)}} \right]$$

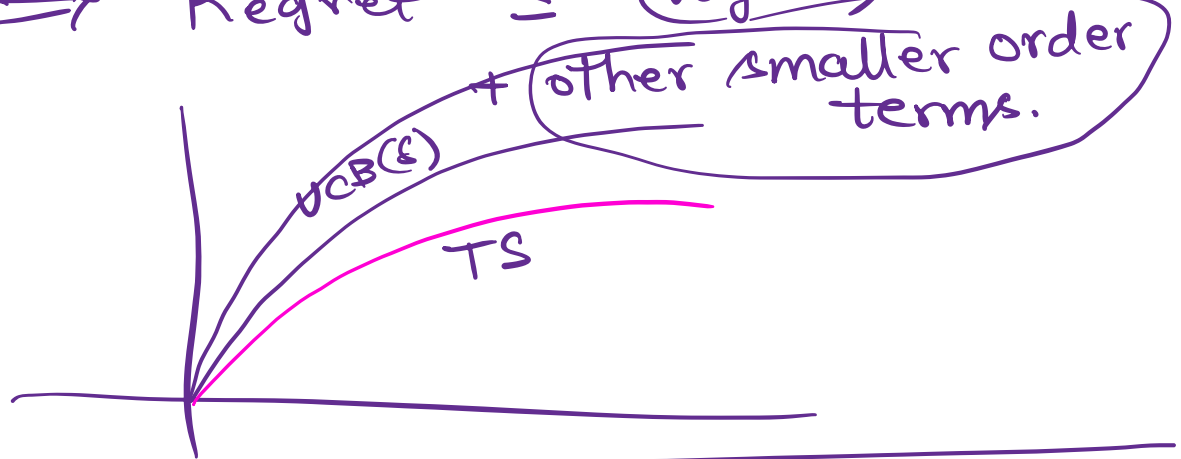
$$0.1 \geq 2 \sqrt{\frac{\log 1/\delta}{2 T_8(t-1)}}$$

$$\Leftrightarrow T_8(t-1) \geq \frac{2 \log 1/\delta}{(0.1)^2}$$

$$\leq \delta$$

$\delta$ : confidence level:  
 IP(picking 0.8 coin instead of 0.9  
 when  $\frac{T}{8} (t-1) \geq \frac{2 \log 1/\delta}{(0.1)^2}$ ]  
 $\leq \delta$

Suppose  $\delta = \frac{1}{t}$   
 Averaged.  
 $\Rightarrow$  Regret  $\leq (\log T) \times \text{constants}$



## General Bandit Set up

$K$  arms

arms = coins

pull = pick

Means:  $\mu_1, \dots, \mu_K$  (you want larger  $\mu$ )

Reward: For arm  $a$ :  
 Reward distribution with mean  $\mu_a$

(e.g. coin case:  
 Reward distribution  
 $= \text{Ber}(p_a)$ )

More generally: Reward diston:  $N(\mu_a, 1)$



$$\text{Reward}(t) = \sum_{s=1}^t \text{Reward in round } s.$$

$$\text{Regret}(t) = t \cdot \max_a \mu_a - \text{Reward}(t)$$

$$\text{Averaged Regret}(t) = \mathbb{E}(\text{Regret}(t))$$

### ① UCB Algorithm:

Assume Reward Distribution has bounded support:  $[a, b]$

Hoeffding (General) to get UCB.

### ② Thompson Sampling



Might be difficult to choose prior. (Natural to try Uniform)

Does not use specific formula for the Reward Distribution.