

## Lecture 2

1. sign-up sheet

2. PTE

3. Intro of me

4. MAB

(1) Recall of RL & MAB setting

RL:  $\mathcal{S}, \mathcal{A}, \mathcal{P}, \Gamma, \gamma$

MAB:  $\mathcal{X}, \mathcal{A}, \mathcal{P}, \Gamma, N$

(\*) time step  $t$ :  $a_k \in \mathcal{A} = \{a_1, \dots, a_K\}$

$$\Gamma_a \sim P(x | \theta_a)$$

we know the reward follows a distribution

e.g.  $\Gamma_{a_i} \sim N(\mu_{a_i}, 1)$

$$\Gamma_{a_i} = \begin{cases} 1 & \text{w.p. } P_i \\ 0 & \text{w.p. } 1 - P_i \end{cases}$$

The optimal arm is  $i^* = \operatorname{argmax}_i \mathbb{E}[\Gamma_{a_i}]$

(e.g.  $i^* = \operatorname{argmax}_i \mu_{a_i}$ ,  $i^* = \operatorname{argmax}_i P_i$ )

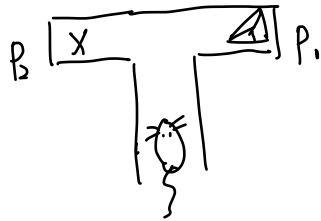
(\*) at the beginning of step  $t$ , we need to make a decision based on the history data observed at

step  $1, \dots, t-1$ .  $\pi_t: \{(A_i, \Gamma_{A_i})\}_{i=1}^{t-1} \rightarrow \mathcal{A}$

(\*) Goal:  $\max \mathbb{E} \left[ \sum_{i=1}^n \Gamma_{A_i} \right]$

## 2) History .

- Introduced by William R. Thomson<sup>P</sup> in 1933.
- Name comes from 1950s by Mosteller & Bush studying animal learning



- Become popular in different applications .
- adaptive experimental design  
(recognized by US Food & Drug Administration)
- new recommendation
- dynamic pricing .
- ad placement → e.g.

### 3). A simple example & a naive approach.

Q: I can put product 1 or 2 at an advertisement place, Not sure the customer likes which. Measure it by click rate. For the next 100 customers who visit the website, what is a good strategy that I can get the most clicks?

why it can be viewed as an MAB prob

$a_1$ : product 1:  $r(a_1) \sim \text{Ber}[p_1]$  ← The prob of clicking the product is  $p_1$   
 $a_2$ : product 2:  $r(a_2) \sim \text{Ber}[p_2]$

- Naive approach:

- Try  $a_1$  25 times

- Try  $a_2$  25 times

- Pull the one with higher empirical means for the rest of 100 times.

- Why it is not optimal?

- when  $\mu_1 = 10$ ,  $\mu_2 = -10$

- when  $\mu_1 = 1$ ,  $\mu_2 = -0.9$

trade off btw exploration & exploitation.

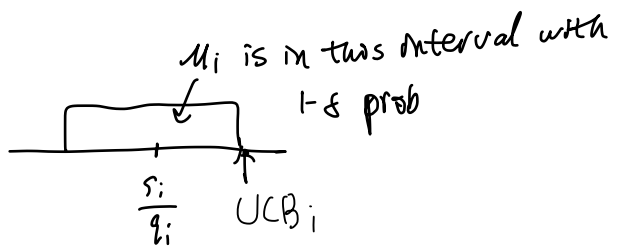
- A better Algorithm:

$$-\frac{S_i}{q_i} + \sqrt{\frac{C}{q_i}}$$

- why it is better?  $\left\{ \begin{array}{l} \text{when } |\mu_1 - \mu_2| \text{ is large} \\ \text{when } |\mu_1 - \mu_2| \text{ is small.} \end{array} \right.$

#### 4). UCB

$$UCB: UCB_i = \frac{S_i}{q_i} + \sqrt{\frac{2 \log(\frac{1}{\delta})}{q_i}}$$

Lemma: For  $\mathcal{I}$ -subgaussian r.v.  $(X - \mathbb{E}X)$ , 

$$\mathbb{P} \left[ \left| \mathbb{E}X - \frac{\sum_{t=1}^i X_t}{i} \right| \leq \sqrt{\frac{2 \log(\frac{1}{\delta})}{i}} \right] \geq 1 - \delta$$

#### Intro of $\sigma$ -subgaussian

##### Formal Def

$\hookrightarrow$   $\sigma$ -subgaussian:  $\iff \mathbb{E}[e^{\lambda X}] \leq e^{\frac{\lambda^2 \sigma^2}{2}}$  for  $\forall \lambda \in \mathbb{R}$

$\hookrightarrow \mathbb{P}(|X| \geq \varepsilon) \leq e^{-\frac{\varepsilon^2}{2\sigma^2}}$  ← The tail decays exponentially fast  
useful ineq. σ is larger, decays faster.

Property of  $\sigma$ -subgaussian:

①  $V[X] \leq \sigma^2$ ,

②  $cX$  is  $|c|\sigma$ -subgaussian

③ If  $X_1, X_2$  are  $\sigma_1, \sigma_2$ -subgaussian

then  $X_1 + X_2$  is  $\sqrt{\sigma_1^2 + \sigma_2^2}$ -subgaussian

Pf of Lemma:  $\sum_{t=1}^i (X_t - \bar{X})$  is  $\sqrt{i}$ -subgaussian

$\frac{1}{i} \sum_{t=1}^i (X_t - \bar{X})$  is  $\frac{1}{\sqrt{i}}$ -subgaussian

$$\begin{aligned} \mathbb{P} \left( \left| \frac{1}{i} \sum_{t=1}^i (X_t - \bar{X}) \right| \geq \sqrt{\frac{2 \log(\frac{1}{\delta})}{i}} \right) &\leq e^{-\frac{2 \log(\frac{1}{\delta})}{\frac{i}{i}}} \\ &= e^{-\log(\frac{1}{\delta})} = \delta \end{aligned}$$

(HW: prove that Normal distribution, Bernoulli distribution, bounded distribution are all  $\sigma$ -subgaussian)

Extension: If  $X - \bar{X}$  is  $\sigma$ -subgaussian, then  $UCB_k(t)$  should

be  $\frac{S_k(t)}{q_k(t)} + \sqrt{\frac{2\sigma^2 \log(\frac{1}{\delta})}{q_k(t)}}$

- What theoretical guarantee do we have?

-  $A = \{a_1, \dots, a_k\}$ .

- reward of  $a_i$  :  $X_{a_i} \sim p(x | \theta_i)$

- Horizon : total # of pulls.

- policy : a mapping from the history data  $\rightarrow$  distribution in action space.  
 $h_{i-1} \rightarrow A_i$

- Goal : maximize  $\mathbb{E} \left[ \sum_{i=1}^n X_{A_i} \right]$

e.g. UCB : history data at the beginning of each round  $i$ ,  
is summarized by  $(s_{t-1}^i, q_{t-1}^i)_{i=1}^k$

$$A_t = \max_i \left\{ \frac{s_{t-1}^i}{q_{t-1}^i} + \sqrt{\frac{2 \log(1/\delta)}{q_{t-1}^i}} \right\}.$$

- regret :  $R_n = n \max_a \mu_a - \mathbb{E} \left[ \sum_{i=1}^n X_{A_i} \right]$

$\Leftrightarrow$  min regret

$\rightarrow$  Thm 2  $k$ -armed  $\delta$ -subgaussian bandit prob.,

for  $\forall$  horizon  $n$ , if  $\delta = \frac{1}{n^2}$ , then

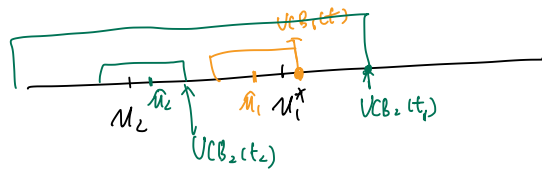
$$R_n = n \max_{a \in A} \mu_a - \mathbb{E} \left[ \sum_{t=1}^n X_t \right] \leq 3 \sum_{i=1}^k \Delta_i + \sum_{i: \Delta_i > 0} \frac{16 \log(n)}{\Delta_i}$$

(where  $\Delta_i = \mu_i - \max_{a \in A} \mu_a$ )

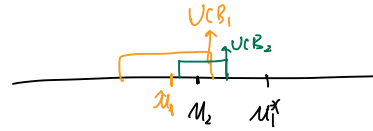
High-level idea: WLOG,  $\mu_1$  is the optimal arm

- When the regret  $> 0$ ? when suboptimal arm  $i > 1$  will be selected

$$UCB_i > \mu_1^*$$



$$UCB_i < \mu_1^*$$



selected at least one of them happens.

1)  $UCB_i(t-1) \geq \mu_1^* \rightarrow$  when this happens sufficiently large,  
 $UCB_i \rightarrow \mu_i < \mu_1^*$  and then it won't happen

2)  $UCB_i(t-1) < \mu_1^* \rightarrow$  this will unlikely happen b/c

$UCB$  is the  $ucb$  of  $i^*$ -th arm.

It happens when  $UCB_i > UCB_1$ .

①  $UCB_i > \mu_1 \leftarrow$  this happens w.p.  $1-\delta$

$$UCB_i > UCB_1 > \mu_1$$

↳ why the # of times this will happen has an upper bd?

↑  
 Subgaussian could give  
 it an explicit upper  
 bd.

$$UCB_i \rightarrow \mu_i < \mu_1 \rightarrow \times$$

②  $UCB_i < \mu_1 \leftarrow$  this happens w.p.  $\frac{\delta}{2}$

"good event"

The UCB value of the optimal arm is always  $> \mu_1$ .

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

$u_i$  to be determined later.

the empirical mean of  $i$ -th arm with  $u_i$  pulls.

⊗ Key pt: exploring  $i$ -th arm  $u_i$  times, its UCB is smaller than the smallest UCB for optimal arm.

↑ After  $u_i$  pulls of  $i$ -th arm, its UCB value  $< \mu_1$ .

We will show two things:

1). If  $G_i$  occurs, then arm  $i$  will be played at most

$u_i$  times:  $T_i(n) \leq u_i$

2). ~~the~~  $G_i^c$  occurs with low prob.  $P(G_i^c) \leq n\delta + e^{-\frac{u_i c^2 \Delta_i^2}{2}}$

holds for  $\forall u_i, c \in (0,1)$

$$\Rightarrow \mathbb{E}[T_i(n)] = \mathbb{E}[T_i(n) | G_i(n)] + \mathbb{E}[T_i(n) | G_i^c(n)]$$

$$\leq u_i + P(G_i^c(n)) n$$

$$\leq u_i + \left( n\delta + e^{-\frac{u_i c^2 \Delta_i^2}{2}} \right)$$

plug in  $u_i = \left\lceil \frac{2 \log(1/\delta)}{(1-c)^2 \Delta_i^2} \right\rceil$  (if  $u_i \geq n$ , ⊗ always holds) &  $\delta = \frac{1}{n^2}$

$$\Rightarrow \mathbb{E}[T_i(n)] = \left\lceil \frac{2 \log(n^2)}{(1-c)^2 \Delta_i^2} \right\rceil + 1 + n^{1-2c^2/(1-c)^2}$$

plug in  $c = \frac{1}{2}$

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

$$R_n = \sum_i \Delta_i \mathbb{E}[T_i(n)] = \sum_i 3 \Delta_i + \frac{16 \log(n)}{\Delta_i}$$



- proof of 1)
- Intro of  $\sigma$ -subgaussian
- proof of 2)
- generalization

pf of 1): Contraction: If  $i$ -th arm is pulled  $u_i$  times, then  $UCB_i < \mu_i$ ,  
 it won't be pulled anymore.

If  $T_i(n) > u_i \Rightarrow \exists t$   
 $\sqrt{T_i(t-1)} = u_i, A_{t-1} = i$

$$\Rightarrow UCB_i(t-1, \delta) = \hat{\mu}_i(t-1) + \sqrt{\frac{2 \log(1/\delta)}{T_i(t-1)}}$$

$$= \hat{\mu}_{i, u_i} + \sqrt{\frac{2 \log(1/\delta)}{u_i}}$$

$$< \mu_i < UCB_i(t-1, \delta)$$

$\Rightarrow A_t = \arg \max_j UCB_j(t-1, \delta) \neq i$  (X)

pf of 2)

$$E_i^c = \underbrace{\{ \mu_i \geq \min_{t \in [n]} UCB_i(t, \delta) \}}_A \cup \underbrace{\{ \hat{\mu}_{i, u_i} + \sqrt{\frac{2 \log(1/\delta)}{u_i}} \geq \mu_i \}}_B$$

$\hookrightarrow$  at least one of  $UCB_i$  value is less than  $\mu_i$

$$\hookrightarrow A \leq \mathbb{P} \left( \bigcup_{s=1}^n \left\{ \mu_i \geq \hat{\mu}_{i, s} + \sqrt{\frac{2 \log(1/\delta)}{s}} \right\} \right)$$

$$\leq \sum_{s=1}^n \mathbb{P} \left( \mu_i \geq \hat{\mu}_{i, s} + \sqrt{\frac{2 \log(1/\delta)}{s}} \right)$$

$$\leq n\delta$$

$\uparrow$   
 when first arm follows  $\sigma$ -subgaussian,  
 this prob  $\leq \delta$

$$B = \mathbb{P} \left( \hat{\mu}_{i:n_i} + \sqrt{\quad} \geq u_i \right)$$

want the empirical mean close to true mean w.h.p

$$= \mathbb{P} \left( \hat{\mu}_{i:n_i} - \mu_i \geq \mu_i - \mu_i - \underbrace{\sqrt{\frac{2 \log(1/\delta)}{u_i}}}_{\substack{\downarrow \\ u_i \text{ large enough s.t.}}}$$

$$\leq \mathbb{P} \left( \hat{\mu}_{i:n_i} - \mu_i \geq c \sigma_i \right) \quad \begin{matrix} < (1-c)(\mu_i - u_i) \\ \uparrow \\ \text{for } c \text{ determined later} \end{matrix}$$

$$\leq e^{-\frac{u_i c^2 \sigma_i^2}{2}} \quad \leftarrow \text{when the } i\text{-th coord follows 1-sub gaussian, the prob has this upper bd.}$$

$$\left( \hat{\mu}_{i:n_i} - \mu_i \text{ is } \frac{1}{\sqrt{u_i}}\text{-sub gaussian} \right)$$

$$\mathbb{P}(G_i^c) \leq n\delta + e^{-\frac{u_i c^2 \sigma_i^2}{2}}$$

Thm 2  $k$ -armed  $\mathcal{I}$ -subgaussian bandit prob.,

for  $\forall$  horizon  $n$ , if  $\delta = \frac{1}{k^2}$ , then

$$R_n = 8\sqrt{kn\log(n)} + 3\sum_i \Delta_i$$

pf:  $R_n = \sum_i \Delta_i \mathbb{E}[T_i(n)]$

$$= \sum_{i: \Delta_i < \Delta} \Delta_i \mathbb{E}[T_i(n)] + \sum_{i: \Delta_i \geq \Delta} 3\Delta_i + \frac{16\log(n)}{\Delta_i}$$

$$= n\Delta + \frac{k16\log(n)}{\Delta} + \sum_{i: \Delta_i \geq \Delta} 3\Delta_i$$

$$\leq 8\sqrt{kn\log(n)} + 3\sum_i \Delta_i$$

UCB Algo

$\delta$  smaller  $\rightarrow$  more exploitation  
 $\delta$  larger  $\rightarrow$  more exploration

$$\text{UCB: } (t-1, \delta) = \left\{ \begin{array}{l} \infty \\ \underbrace{\hat{\mu}_i(t-1)}_{\text{Empirical mean}} + \underbrace{\sqrt{\frac{2\log(t)}{T_i(t-1)}}}_{(1-\delta) \text{ upper bound for } \mathcal{I}\text{-subgaussian}} \end{array} \right.$$

HW: Try UCB on  $N(\mu_1, \sigma)$ ,  $N(\mu_2, \sigma)$  for different  $\Delta = \mu_1 - \mu_2$   
 $\text{Bern}(p_1)$ ,  $\text{Bern}(p_2)$ . for different  $\Delta = p_1 - p_2$ .