

HW3. Pl.2. Piazza @ 29.

online learning w/ bandit feedback:

\mathcal{S} : stochastic multi-armed bandits.

Recall product recommendation:

For $t = 1, 2, \dots, T$:

environment draws $l_t \in [0, 1]^k$ $l_t^{(i)}$:
[not revealed to learner]

$\in [0, 1]^k$ learner selects action $a_t \in \{1, \dots, k\}$.

learner suffers loss $l_t(a_t)$.

$(l(a) : a \in \{1, \dots, k\})$ & sees.
and it does not see $\{l_t(a) : a \neq a_t\}$.

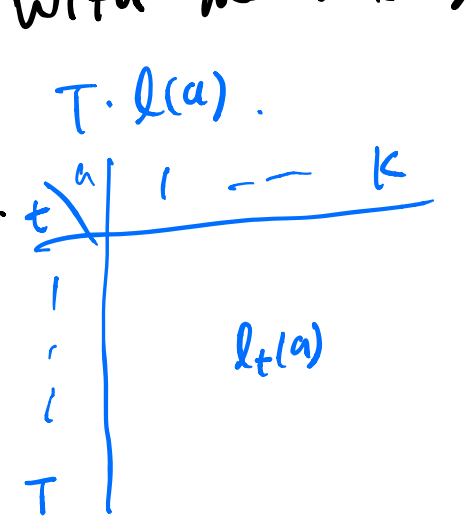
$a \in \mathcal{A}$.

Assume: $\forall a : (l_t(a))_{t=1}^T$ are drawn iid

$\langle A, A \rangle$.

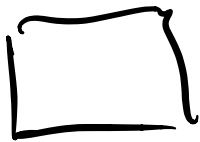
with mean $l(a)$

$l(a) = \sum_{i=1}^K \theta_i a_i$ from a distn over $\{0,1\}^K$.

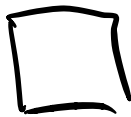


Goal, minimize: $\sum_{t=1}^T l_t(a_t)$

Ex: slot machine game



1



2

...



K

challenge: ① learn which machine is rewarding \rightarrow explore
 ② play machines that are rewarding \rightarrow exploitation

performance metric:

pseudo-regret:

$\mathbb{E} \left[\min_a \sum_{t=1}^T l_t(a) \right]$

$R_T = \mathbb{E} \left[\sum_{t=1}^T l_t(a_t) \right] - T \cdot l(a^*)$

$a^* = \operatorname{argmin}_{a \in \{1, \dots, K\}} l(a)$

$\mathbb{E}_{l_t \sim D} [l_t(a_t)] = l(a_t)$

$\Rightarrow \mathbb{E} [l_t(a_t)] = \mathbb{E} [l(a_t)]$

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

Σ -greedy: w.p. 1 $a_t \sim \text{unif}(\{1, \dots, K\})$
 $1-\epsilon$. $a_t = \arg \min_a \bar{\ell}_t(a)$

$$\bar{\ell}_t(a) = \begin{cases} \frac{\sum_{s=1}^t \mathbb{I}(a_s = a) \ell_s(a_s)}{m_t(a)} & m_t(a) > 0 \\ \odot \text{ I. } & m_t(a) = 0 \end{cases}$$

$$m_t(a) = \sum_{s=1}^t \mathbb{I}(a_s = a) = \# \text{ trials of } a \text{ up to } t$$

w/ appropriate tuning of ϵ . $R_T \leq O(K^{1/3} T^{2/3})$.

Explore-then-commit:

For $t=1, 2, \dots, M$:
 (take action $1, \dots, K$ in round robin) $\approx O(M)$
 $a_t = ((t-1) \bmod K) + 1$

t	a_t
1	1
2	2
...	...
K	K
K+1	1

$\hat{a} = \arg \min_{a \in \{1, \dots, K\}} \bar{\ell}_M(a)$ $\frac{M}{K}$ times
 $\in [\ell(a) \pm \sqrt{\frac{K}{M}}]$

For $t = M+1, \dots, T$:
 $a_t = \hat{a} \rightarrow \sqrt{\frac{K}{M}}$ $\approx O((T-M)\sqrt{\frac{K}{M}})$

roughly speaking $R_T \leq O\left(\frac{M + T\sqrt{\frac{K}{M}}}{\epsilon}\right)$

$\approx O(K^{1/3} T^{2/3})$

$$M = K^3 T^3 \Rightarrow R_T \leq O(K^3 T^3)$$

Analysis:

Lemma: ($K \leq T$), For any online bandit learner, there exists an event E "clean event".

$$P(E) \geq 1 - \frac{2}{T}, \quad e^{[0,1]} \quad e^{[0,1]}$$

$$\forall a, \forall t, \left| \bar{l}_t(a) - l(a) \right| \leq 2 \sqrt{\frac{\ln T}{m_t(a)+1}}$$

handle corner case

DF: uses Hoeffding inequality & union bound ($m_t(a) = 0$).

but is not trivial.

(e.g. Sliuksins: Introduction to MAB)

Regret analysis:

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

$$1 = \underbrace{I(E)}_{\text{clean event}} + I(E^c)$$

$$\leq \underbrace{T \cdot P(E^c)}_{\leq 2} + \mathbb{E} \left[\left(\sum_{t=1}^T l(a_t) - l(a^*) \right) I(E) \right]$$

$$\leq 2.$$

$$\sum_{t=1}^M l(a_t) - l(a^*) + \sum_{t=M+1}^T l(a_t) - l(a^*)$$

$$\leq M \leq (T-M) \cdot 4 \cdot \sqrt{\frac{\ln T}{\frac{M}{K}}}$$

$$\leq 2 + M + 4T \cdot \sqrt{\frac{\ln T \cdot K}{M}}$$

$$M = T^{2/3} K^{1/3}$$

$$= \tilde{O}\left(K^{1/3} T^{2/3}\right)$$

Can we do better?

Idea: using confidence bounds to guide exploration

(Auer, et al, 2002)

UCB algorithm
upper confidence bound

For $t = 1, 2, \dots, T$:

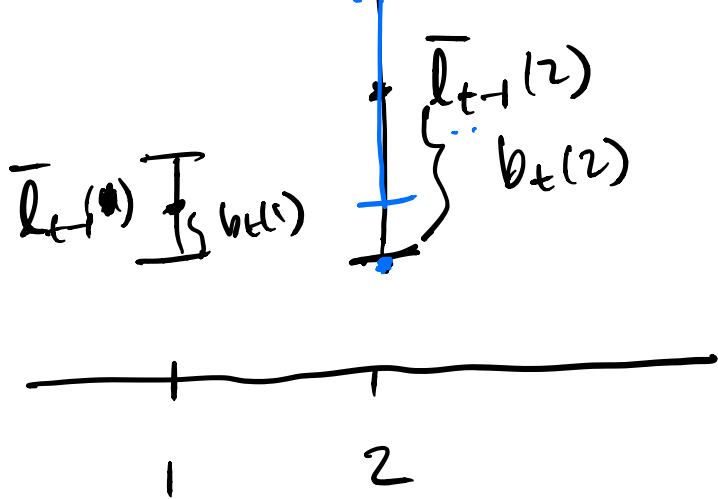
choose $a_t = \underset{a \in \{1, \dots, K\}}{\operatorname{argmin}} \text{LCB}_t(a)$

empirical loss *uncertainty of a*

$$\text{LCB}_t(a) = \underbrace{l_{t-1}(a)} - \underbrace{b_t(a)}$$

$$b_t(a) = 2 \cdot \sqrt{\frac{\ln T}{n_{t-1}(a) + 1}}$$

I



Analysis:

Def, $\Delta(a) = l(a) - l(a^*)$ is the suboptimality gap of action a .

$\Delta(a)$ is large $\Rightarrow a$ is "easy"

$\Delta(a)$ is small \Rightarrow

Thm: UCB satisfies:

$$\textcircled{1} R_T \leq \sum_{a: \Delta(a) > 0} \frac{16 \ln T}{\Delta(a)} + 3K \quad (\text{gap dependent})$$

say $\Delta(a) = \frac{1}{T}$ $\forall a \neq a^*$.

$$\textcircled{2} R_T \leq \mathcal{O}(\sqrt{TK}) \quad (\text{gap independent})$$

Matching lower bound:

Then, \exists const. $c > 0$ for any $\Delta > 0$,
 there exists a stochastic MAB environment,
 such that A satisfies $R_T \geq c \cdot \sqrt{TK}$.

Pf of upper bounds:

Lemma 1: on event \bar{E} :

$\forall a, \forall t,$

(a): $LCB_t(a) \leq l(a)$ (honest)

(b) $LCB_t(a) \geq l(a) - 2b_t(a).$

(tightness)

$$\text{Use } \frac{l(a)}{l(a) - 2b_t(a)}$$

$$\Rightarrow \frac{l_{t+1}(a) - b_{t+1}(a)}{l(a) - 2b_t(a)}$$

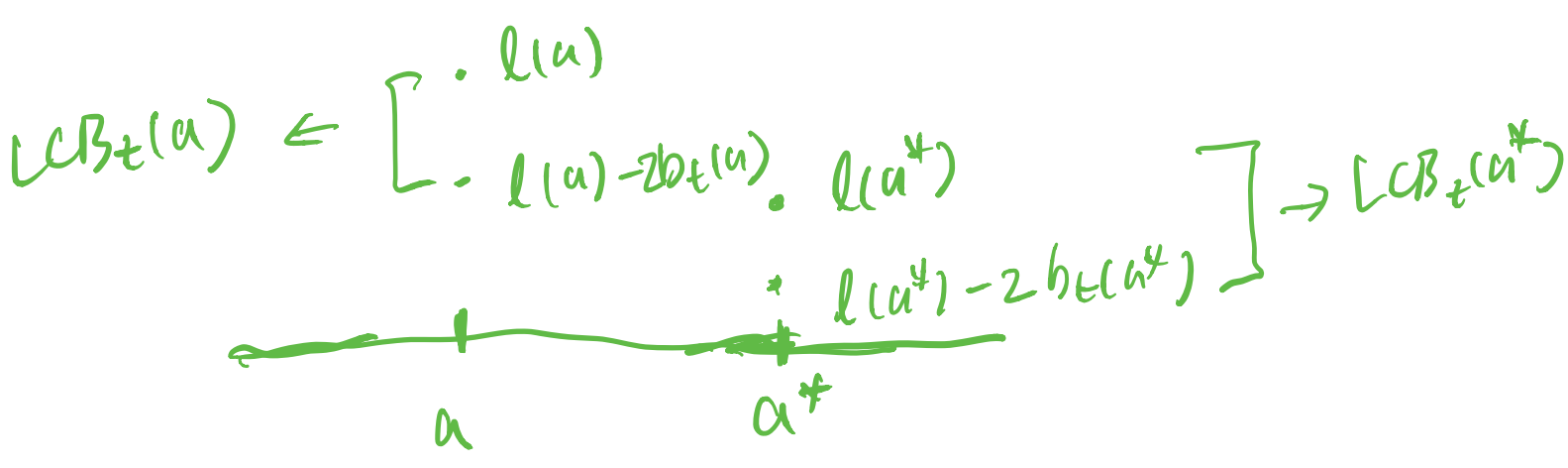
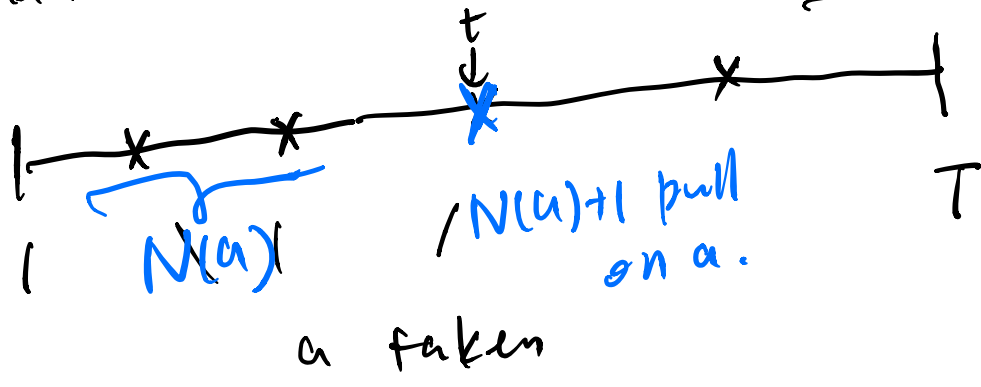
$$\Rightarrow l(a) - b_{t+1}(a) - b_t(a)$$

Lemma 2: $\forall a, \mathbb{E} [m_T(a)] \leq \frac{16 \ln T}{\Delta(a)^2} + 3.$

Pf: key claim: on \bar{E} ,

$$m_T(a) \leq \left\lceil \frac{16 \ln T}{\Delta(a)^2} \right\rceil = N(a).$$

If $m_T(a) \geq N(a) + 1$, then there must exist some t . $a_t = a$, $m_{t+1}(a) = N(a)$.



$m_{t+1}(a) = N(a) \Rightarrow 2b_t(a) < \Delta(a)$.

(exercise).
Lemma 1 (b)

$LCB_t(a) \geq l(a) - 2b_t(a)$

$\geq l(a) - \Delta(a)$

$= l(a^*)$

Lemma 1 (a)
Lemma 1 (a*)

\geq $LCS_t(a)$
 this contradicts w/ $a_t = a$.

$$\Rightarrow \mathbb{E} [m_T(a)]$$

$$= \mathbb{E} [m_T(a) I(E)] + \underbrace{\mathbb{E} [m_T(a) I(E^c)]}_{\leq T \cdot P(E^c)} \leq T$$

$$\leq \left\lceil \frac{16 \ln T}{\Delta(a)^2} \right\rceil + 2$$

$$\lceil x \rceil \leq x + 1.$$

$$\leq \frac{16 \ln T}{\Delta(a)^2} + 3.$$

~~*~~.

concluding the regret analysis:

gap dependent bound:

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

$$= \mathbb{E} \left[\sum_{t=1}^T \Delta(u_t) \right]$$

$$= \mathbb{E} \left[\sum_{a=1}^K \underbrace{\sum_{t: a_t = a} \Delta(u)}_{m_T(a)} \right]$$

$$= \sum_{a=1}^K \Delta(a) \cdot \mathbb{E} [m_T(a)].$$

Lemma 2 $\Rightarrow \sum_{a: \Delta(a) > 0} \Delta(a) \mathbb{E} [m_T(a)]$

$$\leq \sum_{a: \Delta(a) > 0} \Delta(a) \cdot \left(\frac{16 \ln T}{\Delta(a)^2} + 3 \right)$$

$$\Rightarrow \sum_{a: \Delta(a) > 0} \frac{16 \ln T}{\Delta(a)} + 3K.$$

Next time:

- gap in dept bound for UCB
- adversarial MAB problem.