

Noise Free Multi-Armed Bandit Game

Atsuyoshi Nakamura
(Hokkaido University)

David P. Helmbold
(UCSC)

Manfred K. Warmuth
(UCSC)

Noise Free Multi-Armed Bandit Game

for $t = 1, 2, \dots, T$ **do**

- 1 The adversary picks a loss $\ell_{t,i} \in [0, 1]$ for each arm $i \in \{1, \dots, K\}$.
(The player does not know $\ell_{t,i}$ for any arm i .)
- 2 The player chooses arm I_t .
- 3 Loss ℓ_{t,I_t} only is revealed and the player suffers ℓ_{t,I_t} .

Player's obj.: minimizing his/her expected cumulative loss.

Adversary's obj.: maximizing player's expected cumulative loss.

Noise free setting: There is an arm which never suffers any loss.

[Assumptions]

- The player can use a randomized strategy.
- The adversary has infinite computation power.
- The adversary can behave adaptively: the adversary's decision can depend on both the player's past choices and the adversary's past decisions.

- | | |
|-------------------------|---|
| [Assumption] | The losses $\ell_{t,i}$ of all the arms are revealed at the end of time t . |
| [minimax expected loss] | $\sum_{i=2}^K \frac{1}{i} = \Theta(\ln K)$ |
| [Adversary's strategy] | Set the loss of an arm with the highest probability of being chosen to 1 and others to 0. |
| [Player's strategy] | Choose one of the arms with no loss so far randomly with equal probability. |
- Once $\ell_{i,t}$ is set to non-zero and $I_t \neq i$, then the player never chooses i as I_t at time $s > t$.
(In the bandit case, the player may choose it.)

Modification of Noisy-Case Adversary's Strategy

Adversary's strategy Select a best arm i first according to the uniform distribution over $\{1, \dots, K\}$, then generate each component $\ell_{t,j}$ of each loss vector ℓ_t independently according to Bernoulli distribution with parameter 0 for $j = i$ or parameter ϵ for $j \neq i$.

Theorem (Auer et al. 2003)

For the above adversary' strategy, any player algorithm suffers at least

$$\left(1 - \sqrt{\frac{1}{2} \alpha\left(\frac{K}{T}\right)}\right) K - 1$$

loss for $T > \frac{4}{3}K$ in expectation, where $\alpha(x) = -\frac{1}{x} \ln(1 - x)$.

- $\alpha(x)$ is increasing function with $\lim_{x \rightarrow +0} \alpha(x) = 1$.
- The above bound $\rightarrow (1 - 1/\sqrt{2})K - 1$ when $T \rightarrow \infty$.
- There is a player algorithm that suffers at most $(K - 1)/2$ loss for this adversary.

We show an adversary strategy of nose-free multi-armed bandit problem that forces any player algorithm to suffer

$$K - 1 - O(1/T)$$

loss in expectation for K -arm and T -round case.

For any three integers $T' \geq 1$, $K' \geq 2$ and $K' \geq k \geq 1$,

$$F(T', K', k) \stackrel{\text{def}}{=} \left(T' + \binom{K' - 1}{k - 1} - 1 \right) / \left(T' + \binom{K'}{k} - 1 \right).$$

(Example)

$$F(T, 2, 1) = \frac{T}{T + 1}$$

Theorem of Loss Lower Bound

Theorem

For noise free K -armed bandit game, there is an adversary's strategy that forces any player algorithm to suffer the expected loss of at least

$$\sum_{i=1}^m F(t_i - t_{i-1}, K - \sum_{j=1}^{i-1} k_j, k_i)$$

for any positive integers $t_0, \dots, t_m, k_1, \dots, k_m$ that satisfies

$$1 \leq m \leq K - 1, \tag{1}$$

$$1 = t_0 < t_1 < \dots < t_m = T + 1 \text{ and} \tag{2}$$

$$k_1 + \dots + k_m = K - 1. \tag{3}$$

When $K = 2$, uniquely determined as

$$m = 1, t_0 = 1, t_1 = T + 1, k_1 = 1.$$

Corollary

For noise free 2-armed bandit game, there is an adversary's strategy that forces any player algorithm to suffer the expected loss of at least

$$\frac{T}{T+1} = 1 - \frac{1}{T+1}.$$

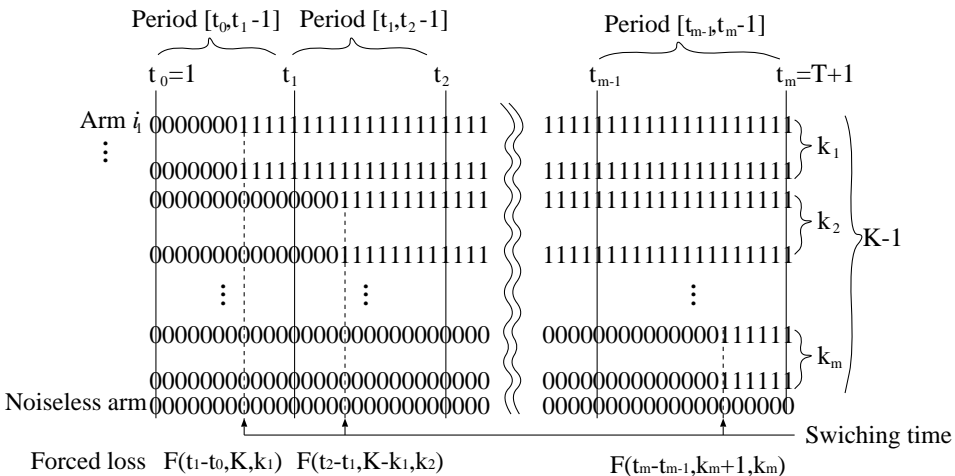
Corollary

For noise free K -armed bandit game, there is an adversary's strategy that forces any player algorithm to suffer the expected loss of at least

$$K - 1 - \frac{K(K - 1)^2(2T + (K - 1)(K + 4))}{(2T + K(K - 1))^2}$$

for $T \geq K(K - 1)/2$.

Adversary's Strategy



Partitioning Problem

Problem

For given two integers $T \geq 1$ and $K \geq 2$, find two non-negative integer sequences t_0, \dots, t_m and k_1, \dots, k_m that maximize

$$\sum_{i=1}^m F(t_i - t_{i-1}, K - \sum_{j=1}^{i-1} k_j, k_i)$$

subject to $1 \leq m \leq K - 1,$ (1)

$1 = t_0 < t_1 < \dots < t_m = T + 1$ and (2)

$k_1 + \dots + k_m = K - 1.$ (3)

Notations

For any natural numbers i, j with $i \leq j$,

$$[i..j] \stackrel{\text{def}}{=} \{i, \dots, j\}, \quad [j] \stackrel{\text{def}}{=} [1..j]$$

For any sequence x_1, \dots, x_n ,

$$\mathbf{x}[b..e] \stackrel{\text{def}}{=} x_b, \dots, x_e.$$

For losses $\ell_{t,i}$ at time t ,

$$\ell_t \stackrel{\text{def}}{=} (\ell_{t,1}, \dots, \ell_{t,K})$$

i_t : realization of random variable I_t

$O_t \stackrel{\text{def}}{=} (I_t, \ell_{t,I_t})$: player's observation at time t

$o_t \stackrel{\text{def}}{=} (i_t, \ell_{t,i_t})$: realization of O_t

For any set $S \subseteq [K]$,

$\mathbb{1}_S \stackrel{\text{def}}{=} K$ -dimensional $\{0, 1\}$ -vector

whose i th component is 1 if and only if $i \in S$.

RepeatW&S[K, T]

Parameter:

K : number of arms

T : number of trials

Initialize:

$t_0, \dots, t_m, k_1, \dots, k_m \leftarrow$ the solution of Partitioning Problem

$d \leftarrow 1$

- 1: **for** time $t = 1, \dots, T$ **do**
- 2: **if** $t \geq t_{d-1}$ **then**
- 3: $b \leftarrow t_{d-1}, e \leftarrow t_d - 1$
- 4: $c \leftarrow \text{BestSwitchingTime}(b, e, \mathbf{o}[1..b - 1], k_d)$
- 5: $d \leftarrow d + 1$
- 6: **end if**
- 7: $\ell_t \leftarrow \text{Wait\&Sticking}[b, c, e, \mathbf{o}[1..b - 1], k_d](t, \mathbf{o}[b, ..t - 1])$
- 8: Observe the player's choice i_t
- 9: **end for**

BestSwitchingTime($b, e, \mathbf{o}[1..b-1], k$)

Input: $b, e \in \mathbb{N}$: beginning and ending times with $1 \leq b \leq e$
 $\mathbf{o}[1..b-1]$: players observations from time 1 to time $b-1$
 k : number of no-loss arms to switch to 1-loss

Output: c^* : best time to switch from waiting to sticking

1: $S \leftarrow$ the set of arms with no loss by time $b-1$

2: $p_{\max} = -1$

3: **for** $c = b, \dots, e$ **do**

4: $p_c \leftarrow$

$$E_{I[b..e]} \left[\max_{s \subseteq S, |s|=k} P \left\{ \sum_{t=c}^e \ell_{t, I_t} \geq 1 \left| \begin{array}{l} \mathbf{O}[1..b-1] = \mathbf{o}[1..b-1], \\ I[b..c-1], \\ \ell_b = \dots = \ell_{c-1} = \mathbb{1}_{[K] \setminus S}, \\ \ell_c = \dots = \ell_e = \mathbb{1}_{([K] \setminus S) \cup s} \end{array} \right. \right\} \right]$$

5: **if** $p_c > p_{\max}$ **then**

6: $c^* \leftarrow c, p_{\max} \leftarrow p_c$

7: **end if**

8: **end for**

9: **return** c^*

Lemma

Let b and T' be arbitrary positive integers and let k be an arbitrary non-negative integer. Let c^* be the returned value from $\text{BestSwitchingTime}(b, b + T' - 1, \mathbf{o}[1..b - 1], k)$. Then, for any player algorithm, the following holds with respect to the loss vectors generated by

$\text{Wait\&Sticking}[b, c^*, b + T' - 1, \mathbf{o}[1..b - 1]](t, \mathbf{o}[b..t - 1])$:

$$\mathbb{E}_{I[b..b+T'-1]} \left[\sum_{t=b}^{b+T'-1} \ell_{t, I_t} \right] \geq F(T', |S|, k),$$

where S is the set of arms with no loss by time $b - 1$, that is, $S = \{i \in [K] : \sum_{t=1}^{b-1} \ell_{t,i} = 0\}$.

Proof Scketch (1/2)

BestSwitchingTime($b, e, \mathbf{o}[1..b-1], k$) returns $c^* = c$ that maximizes

$$p_c \stackrel{\text{def}}{=} E_{I[b..e]} \left[\max_{s \subseteq S, |s|=k} P \left\{ \sum_{t=c}^e \ell_{t, I_t} \geq 1 \left| \begin{array}{l} \mathbf{O}[1..b-1] = \mathbf{o}[1..b-1], \\ I[b..c-1], \\ \ell_b = \dots = \ell_{c-1} = \mathbb{1}_{[K] \setminus S}, \\ \ell_c = \dots = \ell_e = \mathbb{1}_{([K] \setminus S) \cup s} \end{array} \right. \right\} \right]$$

Define

$$p_{b,s}(\mathbf{o}[1..b-1], e-b+1) \stackrel{\text{def}}{=} P \left\{ \sum_{t=b}^e \ell_{t, I_t} \geq 1 \mid \begin{array}{l} \mathbf{O}[1..b-1] = \mathbf{o}[1..b-1], \\ \ell_b = \dots = \ell_e = \mathbb{1}_{([K] \setminus S) \cup s} \end{array} \right\}.$$

Then,

$$p_b = \max_{s \subseteq S, |s|=k} p_{b,s}(\mathbf{o}[1..b-1], e-b+1).$$

Proof Scketch (2/2)

Thus,

$$\begin{aligned} \mathbb{E}_{I[b..e]} \left[\sum_{t=b}^e \ell_{t, I_t} \right] &\geq p_{c^*} \\ &= \max \left(\{p_{b,s}(\mathbf{o}[1..b-1], e-b+1) \mid s \subseteq S, |s| = k\} \cup \{p_c \mid c = b+1, \dots, e\} \right) \end{aligned}$$

We show

$$\sum_{s \subseteq S, |s|=k} p_{b,s}(\mathbf{o}[1..b-1], e-b+1) + \sum_{c=b+1}^e p_c \geq e-b+1 + \binom{|S|-1}{k-1} - 1,$$

which implies

$$p_{c^*} \geq \left(e-b + \binom{|S|-1}{k-1} \right) / \left(e-b + \binom{|S|}{k} \right)$$

What is the minimax-strategy of the noise-free multi-armed bandit game?

Minimax Strategy for Noise Free 2-Armed Bandit Game

```
1:  $i^* \leftarrow ?$ 
2: for time  $t = 1, \dots, T$  do
3:   if  $i^* = ?$  then
4:     if  $t > 1$  and  $i_{t-1} = 1$  then
5:        $I_t \leftarrow 1$ 
6:     else
7:        $i_t \leftarrow \begin{cases} 1 & \text{with prob. } 1/(T-t+2) \\ 2 & \text{with prob. } (T-t+1)/(T-t+2) \end{cases}$ 
8:     end if
9:   else
10:     $i_t \leftarrow i^*$ 
11:  end if
12:   $l_{t,i_t}$  is observed.
13:  if  $l_{t,i_t} > 0$  then
14:     $i^* \leftarrow \begin{cases} 1 & (i_t = 2) \\ 2 & (i_t = 1) \end{cases}$ 
15:  end if
16: end for
```