

## Bandits with Oblivious Experts.

Source: Bandit Algorithms, L&S, Chap 18.2.

Setting: Time horizon =  $n$

$k$ -armed bandit, adversarial environment

$M$  experts

an expert recommends arms, i.e.,

For expert  $m$ :  $(e_{m1}, e_{m2}, \dots, e_{mk})$  is  
a probability distribution over the  $k$  arms,  
with  $0 \leq e_{mi} \leq 1$ ,  $\sum_{i=1}^k e_{mi} = 1$ .

At time  $t$ :

$$E^{(t)} = \begin{pmatrix} e_{11}^{(t)} & \cdots & \cdots & \cdots & e_{1k}^{(t)} \\ e_{21}^{(t)} & e_{22}^{(t)} & \cdots & \cdots & e_{2k}^{(t)} \\ \vdots & & & & \vdots \\ e_{m1}^{(t)} & e_{m2}^{(t)} & \cdots & \cdots & e_{mk}^{(t)} \\ \vdots & & & & \vdots \\ e_{M1}^{(t)} & \cdots & \cdots & \cdots & e_{Mk}^{(t)} \end{pmatrix}$$

$\leftarrow k \text{ arms} \rightarrow$

↑  
 $M$   
experts  
↓

Oblivious Expert: Experts recommendations can  
vary with time, BUT experts cannot learn

from past observations of this bandit player.

e.g.: Experts are pre-trained and static:  
 $E^{(t)} = E$ .

### Player's Selection of Action:

Player chooses a distribution over Experts, i.e.,

\*  $Q_t = (q_1^{(t)}, \dots, q_M^{(t)})$ , a prob. dist.  
over  $\{1, 2, \dots, M\}$  experts.

\* Expert  $m$  recommends arm  $j$  from  $m^{\text{th}}$  row  
of  $E^{(t)}$  ( $E^{(t)}$  known to player).

\* Effectively, player plays dist  $P_t$  over arms  
 $\{1, 2, \dots, k\}$ , where

$$P_t = Q_t \cdot E^{(t)} = (p_1^{(t)}, \dots, p_k^{(t)}).$$

Player chooses policy  $\pi = (\pi_1, \pi_2, \dots, \pi_n)$

where  $\pi_t : \text{History}(t-1) \mapsto \{1, \dots, M\}$ ,

i.e., at time  $t$ , chooses a dist. over experts  
(which in turn, induces a dist. over arms).

Adversarial Setting:

action at time 1  
reward at time 1  $\in [0,1]$ .

- ① At beginning of time  $t$ :  $(A_1, X_1, \dots, A_{t-1}, X_{t-1})$   
known to adversary, player.

Player has policy  $\pi$  (known to player and adversary)

- ② Adversary chooses arm rewards  $(x_t, \dots, x_k)$ ,  
with  $x_{tj} \in [0,1]$  (or equivalently, chooses  
loss  $y_{tj} = 1 - x_{tj}$ ).

- ③ Player uses (secret) randomness to sample an  
expert according to dist.  $Q_t$ , and further uses  
(secret) randomness to sample an arm from the  
chosen expert's dist. Equivalently, player samples  
arm  $j$  from  $P_t$ ,

$$\text{i.e. } A_t \sim P_t = Q_t E^{(t)}$$

④ Player observes reward  $X_t = x_{tA_t}$

$$\text{Regret: } R_n(\pi, x) = \max_{1 \leq m \leq M} \sum_{t=1}^n E_m^{(t)} x_t - \mathbb{E} \left[ \sum_{t=1}^n x_t \right],$$

i.e., with respect to the best fixed expert in hindsight.

Exp 4: Exponential Weighting for Exploration and Exploitation with Experts.

Input:  $n, k, M, \eta$  (see text for a slightly more general version).

Initialize: Choose experts equally likely, i.e.,

$$Q_1 = \left( \frac{1}{M} \quad \frac{1}{M} \quad \cdots \quad \frac{1}{M} \right).$$

for each  $t \geq 1, t = 1, 2, \dots, n$ :

⑤ Experts reveal  $E^{(t)}$  (to player and adversary)

(b) Player chooses arm  $A_t \sim Q_t E^{(t)}$

(c) Player sees reward  $X_t = x_{tA_t}$

(d) Importance Sampling:  $\hat{\gamma}_{tj} = \left( \frac{x_{\{A_t=j\}}}{P_{tj}} \right) Y_t$ ,

where  $Y_t = (1 - X_t)$  and  $\hat{x}_{tj} = (1 - \hat{\gamma}_{tj})$

(e)  $\tilde{x}_t = E^{(t)} \hat{x}_t$  ( $\xrightarrow{k\text{-dim. vector}}$   
 $\hookrightarrow M\text{-dim. vector}$ ) (expert-level estimate of rewards).

(f) Update  $Q_t$ :  $Q_{t+1,m} = \left( \frac{e^{n \tilde{x}_{tm}} \cdot Q_{tm}}{\sum_l e^{n \tilde{x}_{tl}} \cdot Q_{tl}} \right)$ ,

$$m = 1, 2, \dots, M$$

Theorem (18.1 in textbook):  $\eta = \sqrt{\frac{2 \ln(M)}{nk}}$ . Then,

$$R_n(\pi, x) \leq \sqrt{2nk \ln(M)}$$

Proof: We use the following corollary (that is from the proof of Exp 3, see Thm 11.2 in textbook). This is a sharpened version that we saw at the end of Step B in the proof of Exp 3.

-----  
Lemma (18.2 in text): Fix any expert  $m^* \in \{1, 2, \dots, M\}$ . Then,

$$\left( \sum_{t=1}^n \tilde{x}_{tm^*} - \sum_{t=1}^n \sum_{m=1}^M Q_{tm} \tilde{x}_{tm} \right)$$

$$\leq \frac{\ln(M)}{\eta} + \frac{n}{2} \sum_{t=1}^n \sum_{m=1}^M Q_{tm} (1 - \tilde{x}_{tm})^2$$

the improvement by a factor of 2 follows from a tighter analog of Step B in Exp 3, available in Thm 11.2 in textbook.

At the begining of time  $t$ , history is :

$$I_t = \{E^{(1)}, A_1, X_1, E^{(2)}, A_2, X_2, \dots, E^{(t-1)}, A_{t-1}, X_{t-1}, E^{(t)}\}$$

Let  $m^*$  be the best arm in hindsight; then,

$$m^* = \arg \max_{1 \leq m \leq M} \sum_{t=1}^n E_m^{(t)} x_t$$

(  $x_{t1}$   
⋮  
 $x_{tk}$  ) = arm rewards at time t  
m<sup>th</sup> row of  $E^{(t)}$

From Lemma above,

$$\begin{aligned} \sum_{t=1}^n \tilde{x}_{tm^*} - \sum_{t=1}^n \sum_{m=1}^M Q_{tm} \tilde{x}_{tm} \\ \leq \frac{\ln(M)}{\eta} + \frac{1}{2} \sum_{t=1}^n \sum_{m=1}^M (Q_{tm}(1 - \tilde{x}_{tm}))^2 \end{aligned}$$

L  $\rightarrow$  (★).

Recall that  $\hat{x}_{ti}$  is unbiased, i.e.,

$$\mathbb{E}[\hat{x}_{ti} | I_t] = x_{ti}$$

$$\Rightarrow \mathbb{E}[\tilde{x}_t | I_t] = \mathbb{E}[E^{(t)} \hat{x}_t | I_t]$$

$\underbrace{\text{M-dim. column vector.}}_{\text{(Mxk) matrix}} \quad \underbrace{\text{k-dim. column vector}}_{\text{unbiased}}$

$$= E^{(t)} \underbrace{\mathbb{E}[\hat{x}_t | I_t]}_{\text{unbiased}} = E^{(t)} \cdot x_t$$

Taking  $\mathbb{E}[\cdot]$  on both sides of  $\textcircled{*}$ :

$$R_n(\pi, x) \leq \frac{\ln(M)}{\gamma} + \frac{\eta}{2} \sum_{t=1}^n \sum_{m=1}^M \mathbb{E} \left[ Q_{tm}(r - \tilde{x}_{tm})^2 \right]$$

Recall:  $\hat{\gamma}_{ti} = (1 - \hat{x}_{ti})$ ,  $y_{ti} = (1 - x_{ti})$ ,  $\tilde{\gamma}_{tm} = (1 - \tilde{x}_{tm})$

time  $\swarrow$  arm  $\searrow$    
 $\hat{\gamma}_{ti}$  loss reward.   
 arm  $\swarrow$  time  $\searrow$  expert  $\tilde{\gamma}_{tm}$

Also,  $\tilde{\gamma}_t = E^{(t)} \hat{\gamma}_t$ ,  $\hat{\gamma}_{ti} = \left( \frac{\sum_{\{A_t=i\}} y_{ti}}{P_{ti}} \right)$

$$\therefore \mathbb{E} \left[ \tilde{\gamma}_{tm}^2 \mid I_t \right] = \mathbb{E} \left[ \left( \frac{E_{mA_t}^{(t)} y_{tA_t}}{P_{tA_t}} \right)^2 \mid I_t \right]$$

$$= \sum_{j=1}^k \left( \dots \right) \cdot P_{tj} = \sum_{j=1}^k \left( \frac{E_{mj}^{(t)} y_{tj}}{P_{tj}} \right)^2$$

$P(A_t=j)$

$$\leq \sum_{j=1}^k \frac{E_{mj}^{(t)}}{P_{tj}} \quad \left( \because 0 \leq E_{mj}^{(t)}, y_{tj} \leq 1 \right)$$

$$\Rightarrow \mathbb{E} \left[ \sum_{m=1}^M Q_{tm} (1 - \tilde{x}_{tm})^2 \right] = \mathbb{E} \left[ \sum_{m=1}^M Q_{tm} \tilde{y}_{tm}^2 \right]$$

$$\leq \mathbb{E} \left[ \sum_{m=1}^M Q_{tm} \mathbb{E} \left[ \tilde{y}_{tm}^2 \mid I_t \right] \right]$$

$$\leq \mathbb{E} \left[ \sum_{m=1}^M Q_{tm} \cdot \sum_{j=1}^K \frac{E_{mj}^{(t)}}{P_{tj}} \right]$$

$$= \mathbb{E} \left[ \sum_{j=1}^K \underbrace{\frac{\sum_{m=1}^M Q_{tm} E_{mj}^{(t)}}{P_{tj}}}_{P_{tj}} \right]$$

$$= K$$

$$\Rightarrow R_n(\pi, x) \leq \frac{\ln(M)}{\eta} + \frac{\eta n k}{2}$$

$$\text{Choose } \eta = \sqrt{\frac{\ln(M)}{nk}} \Rightarrow R_n(\pi, x) \leq \sqrt{2nk \ln(M)}$$

□