

EE6106

## Lecture 4

(Date: 23 January 2024)

### Adversarial Bandits

The information table can be written as an array of the form

$$y = \begin{bmatrix} y_{1,1} & \dots & y_{1,T} \\ \vdots & y_{i,T} & \\ y_{k,1} & \dots & y_{k,T} \end{bmatrix}$$

We are not assuming  
any kind of underlying  
distribution for the losses

loss of expert / arm  $i$   
at time  $T$

Randomization is necessary in the presence of an adversary.

If the environment is however, structured we can have deterministic algo as well.

### Protocol

1. At time  $t \geq 1$ ,

- algo picks arm  $A_t \in [K]$
- loss  $Y_{A_t, t}$  is revealed

### Goal

Minimize  $\mathbb{E} \left[ \sum_{t=1}^T Y_{A_t, t} \right] - \left( \min_i \sum_{t=1}^T Y_{i, t} \right)$  ← try and make this sublinear

let's say start calling this  
to be " $\mathbb{E} R_T$ " instead of  $\mathbb{E} CRT$

Idea : Replace  $Y_{j, t}$  by estimate  $\hat{Y}_{j, t}$ , which is an unbiased estimator  
(conditionally unbiased)

Say  $p_{j, t}$  = probability of picking arm  $j$  at time  $t$  given  $I_{t-1}$

$$I_{t-1} = (A_1, Y_1, A_2, Y_2, \dots, A_{t-1}, Y_{t-1})$$

$Y_t$  : observed loss, and hence  $Y_t = Y_{A_t, t}$

$$\hat{Y}_{j, t} = \frac{\mathbb{1}_{\{A_t = j\}} Y_t}{p_{j, t}}$$

indicator of subsequently actually picking action  $j$

$$= \begin{cases} 0 & i \neq A_t \\ \frac{Y_{i, t}}{p_{j, t}} & j = A_t \end{cases}$$

We know that  $\mathbb{E}[\hat{Y}_{j, t} | I_{t-1}] = Y_{j, t}$

The algorithm does not always have access to  $Y_{jt}$  so we can instead use  $\hat{Y}_{jt}$ .

$$\text{Var}(\hat{Y}_{jt} | I_{t-1}) = \frac{Y_{jt}^2 (1 - p_{jt})}{p_{jt}}$$

If  $p_{jt} \rightarrow 0$ , then the variance tends to be QUITE LARGE

$$\hat{L}_{jt} = \sum_1^T \hat{Y}_{jt}$$

### EXP3 (EXPONENTIAL ALGORITHM FOR EXPLORE & EXPLOIT)

(There is also an EXP4 algorithm, lol)

At time  $t \geq 1$ ,

$$p_{jt} = \frac{e^{-\eta \hat{L}_{jt,t-1}}}{\sum_i e^{-\eta \hat{L}_{i,t-1}}}$$

Pick arm  $A_t \sim (p_{jt})_{j=1}^K$

$$\hat{L}_{jt} = \hat{L}_{j,t-1} + \hat{Y}_{jt} \quad \dots \text{(update step)}$$

If we replace  $\hat{Y}$  by  $Y$ , then this algorithm is exactly similar to REWMA

THEOREM For Exp 3,

$$R_T \leq \frac{\log K}{\eta} + \eta TK$$

with  $\eta = \sqrt{\frac{\log K}{TK}}$ ,  $R_T \leq 2\sqrt{TK \log K}$

PROOF :  $R_{i,T} = \mathbb{E} \left[ \sum_{t=1}^T Y_t \right] - \sum_i Y_{i,t}$  Expected loss of algorithm relative to an arbitrarily selected arm

Step 1  $\hat{L}_{i,t} = \sum_{s=1}^t \hat{Y}_{i,s} \rightarrow$  estimate of loss of algorithm ~~up to~~ <sup>of</sup> arm  $i$

$$\mathbb{E}[\hat{L}_{i,t}] = \sum_{s=1}^T \mathbb{E}[\mathbb{E}[\hat{Y}_{i,s} | I_{s-1}]]$$



$$\therefore \mathbb{E}[L_{i,T}] = \sum_{t=1}^T Y_{i,t} \quad \dots \quad (1)$$

$$\mathbb{E}[Y_t | I_{t-1}] = \sum_{j=1}^K p_{j,t} Y_{j,t}$$

$$\hat{L}_T = \sum_{t=1}^T \sum_{j=1}^K p_{j,t} \hat{Y}_{j,t}$$

$$\mathbb{E}[\hat{L}_T] = \sum_{t=1}^T \sum_{j=1}^K \mathbb{E}[\mathbb{E}[p_{j,t} \hat{Y}_{j,t} | I_{t-1}]]$$

$$= \sum_{t=1}^T \sum_{j=1}^K \mathbb{E}[p_{j,t} Y_{j,t}]$$

$$= \sum_{t=1}^T \mathbb{E}[Y_t]$$

We shall then write

$$R_{i,T} = \mathbb{E}[\hat{L}_T] - \mathbb{E}[L_{i,T}]$$

→ regret of algo rel to arm is NOT regret of arm.

Step 2

$$W_t = e^{\eta t} \sum_{j=1}^K e^{-\eta \hat{L}_{j,t}}$$

Lower Bound :  $W_T \geq e^{\eta T} e^{-\eta \hat{L}_{i,T}}$

Upper Bound :  $\frac{W_t}{W_{t-1}} = e^{\eta} \left( \sum_{j=1}^K e^{-\eta \hat{L}_{j,t-1}} \right) e^{-\eta \hat{Y}_{i,t}}$

$$\sum_{j=1}^K e^{-\eta \hat{L}_{j,t-1}}$$

→ probability of choosing arm

$$= e^{-\eta} \sum_{j=1}^K p_{j,t} e^{-\eta \hat{Y}_{j,t}}$$

$$= \sum_{j=1}^K p_{j,t} e^{\eta(1 - \hat{Y}_{j,t})}$$

For  $x \leq 1$ ,

$$1 + x \leq e^x \leq 1 + x + x^2$$

$$\therefore \frac{W_t}{W_{t-1}} \leq \sum p_{jt} (1 + \eta(1 - \hat{y}_{jt}) + \eta^2(1 - \hat{y}_{jt})^2)$$

$$= 1 + \eta \sum p_{jt} (1 - \hat{y}_{jt}) + \eta^2 \sum p_{jt} (1 - \hat{y}_{jt})^2$$

$$\frac{W_t}{W_{t-1}} \leq e^{\eta \sum p_{jt} (1 - \hat{y}_{jt}) + \eta^2 \sum p_{jt} (1 - \hat{y}_{jt})^2}$$

summation is wrt

the index of arms  $\{1, \dots, k\}$

$$\frac{W_T}{K} = \prod_{t=1}^T \frac{W_t}{W_{t-1}} \leq e^{\eta \sum \sum p_{jt} (1 - \hat{y}_{jt}) + \eta^2 \sum \sum p_{jt} (1 - \hat{y}_{jt})^2}$$

call this  $N$

Combining the LB & UB, we have

$$e^{\eta T} e^{-\eta L_{i,T}} \leq K N \quad \begin{matrix} \sum \sum p_{jt} & \sum \sum p_{jt} \hat{y}_{jt} \\ \uparrow & \uparrow \\ \sum \sum p_{jt} (1 - \hat{y}_{jt}) & \sum \sum p_{jt} \hat{y}_{jt} \end{matrix}$$

Taking logarithms

$$\eta T - \eta L_{i,T} \leq \log K + \eta (T - \hat{L}_T) + \eta^2 \sum \sum p_{jt} (1 - \hat{y}_{jt})^2$$

$$\therefore T - L_{i,T} \leq \frac{\log K}{\eta} + T - \hat{L}_T + \eta \sum_{t=1}^T \sum_{j=1}^k p_{jt} (1 - \hat{y}_{jt})^2$$

Note:  $\hat{y}_{jt}$  can be arbitrarily large though  $y_{jt} \in [0, 1]$   
 we wish to bound the expectation on the RHS

Lemma:  $\mathbb{E} \left[ \sum \sum p_{jt} (1 - \hat{y}_{jt})^2 \right] \leq KT$

$$\therefore \hat{L}_T - L_{i,T} \leq \frac{\log K}{\eta} + \eta \sum_{t=1}^T \sum_{j=1}^k p_{jt} (1 - \hat{y}_{jt})^2$$

Proof of lemma:

$$\mathbb{E} \left[ \sum \sum p_{jt} (1 - \hat{y}_{jt})^2 \right] = \mathbb{E} \left[ \sum_{t=1}^T \sum_{j=1}^k p_{jt} (1 + \hat{y}_{jt}^2 - 2\hat{y}_{jt}) \right]$$



$$\begin{aligned}
 \mathbb{E} \left[ \sum_t \sum_j p_{jt} (1 - \hat{y}_{jt})^2 \right] &= \sum_{t=1}^T \mathbb{E} \left[ 1 - 2Y_t + \sum_j p_{jt} \mathbb{1}_{\{A_t=j\}} Y_{jt}^2 \right] \\
 &= \sum_{t=1}^T \mathbb{E} \left[ 1 - 2Y_t + \sum_j Y_{jt}^2 \right] \quad \text{because } \mathbb{E} \left[ \sum_j p_{jt} \mathbb{1}_{\{A_t=j\}} Y_{jt}^2 \right] = \mathbb{E} [Y_t^2] \\
 &= \sum_{t=1}^T \mathbb{E} \left[ (1 - Y_t)^2 + \sum_{j \neq A_t} Y_{jt}^2 \right] \\
 &\quad \underbrace{\leq 1}_{\text{Just completing squares}} \quad \underbrace{\leq (k-1)}_{\text{random since } A_t \text{ is random}} \\
 &\leq kT
 \end{aligned}$$

Hence  $R_T \leq \frac{\log K}{\eta} + \eta TK$ .

Note: A subtle difference in the bounds between REWMA & EXP3 is that there is a factor of  $K$  appearing here. Thus, the regret is of the order  $O(\sqrt{k})$  <sup>worse</sup> as compared to the full information setting.

choose  $\eta = \sqrt{\frac{\log K}{TK}}$ , we get  $\rightarrow R_T \leq 2\sqrt{TK \log K}$

The actual payment comes from here

①  $\mathbb{E} \left[ \sum_t \sum_j p_{jt} (1 - \hat{y}_{jt})^2 \right] \leq kT$

② However, if we replace  $\hat{y}_{jt}$  by  $y_{jt}$ , so as to get to the full information setting, the bound we would have got is

$$\begin{aligned}
 \mathbb{E} \left[ \sum_t \sum_j p_{jt} (1 - y_{jt})^2 \right] &\leq T \\
 &\downarrow \\
 &\leq 1
 \end{aligned}$$

Note: If the losses are unbounded, the nature of bounds obtained are much much weaker

Ref: Lattimore presents identical proofs (From now onwards shall be following Lattimore)

M	T	W	T	F	S	S
Page No.:						YOUVA
Date:						

The high probability bounds for the regrets obtained in this algorithm are quite bad. There exists a refined algorithm for this called the  $\text{Exp3-IX}$