

## ONLINE LEARNING / ADAPTIVE DECISION MAKING

Set of actions  $\mathcal{X}$  (e.g.  $\mathcal{X} = \{1, \dots, k\}$ ,  $\mathcal{X} \subseteq \mathbb{R}^d$ )

For  $t = 1 : T$

Adversary picks loss  $l_t \in \mathcal{F} : \mathcal{X} \rightarrow \mathbb{R}$  unknown  
 DM picks  $x_t \in \mathcal{X}$   
 DM observes something about  $l_t$

Goal:  $\sum_{t=1}^T l_t(x_t) \rightarrow \min!$

What does DM observe?

Full information setting:  $l_t$

Bandit setting:  $l_t(x_t)$

Performance metric

Regret  $R_T = \sum_{t=1}^T l_t(x_t) - \min_{x \in \mathcal{X}} l_t(x)$

Want  $R_T/T \rightarrow 0$

"average excess loss compared to best fixed action in hindsight vanishes!"

$\mathcal{E}$ :  $\mathcal{X}$  = set of experts  
classifiers

$h_1, \dots, h_E$

At each round, receive data point, and predictions  $y_t(1), \dots, y_t(E)$  for the experts  
 $y_t(x) \in \{0, 1\}$

We predict  $y_t$ , then observe  $y_t^*$

$l_t(x) = \begin{cases} 0 & \text{if } y_t(x) = y_t^* \\ 1 & \text{otherwise} \end{cases}$

"full information" setting

## FULL INFORMATION

Warmup: Assume  $X = \{1..k\}$   
 $l_t \in \{0, 1\}$   $b_{x,t}$   
 $\exists x: l_t(x) = 0 \quad \forall t$

Holiving alg.: maintain weights per expert

$$w_t(x) = 1 \quad \forall x \in X$$

At time  $t$ : let

$$\forall y \in \{0, 1\}: w_t(y) := \sum_{x: y_t(x) = y} w_t(x)$$

total weight  
of all experts  
voting for  $y$

Predict weighted majority

$$\hat{y}_t := \operatorname{argmax}_{y \in \{0, 1\}} w_t(y)$$

Pick some  $x_t$  who predicted  $\hat{y}_t$

Receive  $y_t^*$

$$\forall x: \text{set } w_{t+1}(x) = \begin{cases} 0 & \text{if } \hat{y}_t \neq y_t^* \\ w_t(x) & \end{cases}$$

Claim:  $R_T \leq \lceil \log_2 k \rceil$

Proof: Each round either no mistake ( $l_t(x_t) = 0$ )

$$\text{OR } \sum_{x=1}^k w_{t+1}(x) \leq \frac{1}{2} \sum_{x=1}^k w_t(x)$$

□

More generally

What if no expert is always correct?

What if  $l_t(x) \in [0, 1]$  ?

Multiplicative weights / Hedge:

$$\text{Init } w_1(x) = 1 \quad \forall x$$

for  $t = 1 \dots T$

$$p_t(x) = \frac{w_t(x)}{\sum_{x'} w_t(x')}$$

pick  $x_t \sim p_t$

Incur loss  $l_t(x_t)$

$$w_{t+1}(x) = w_t(x) \cdot \exp(-\varepsilon l_t(x)) \quad \forall x$$

If  $l_t \in [0, 1]$   
for "incorrect"  
experts, reduce  
weight  
 $\exp(-\varepsilon)$   
 $\approx (-\varepsilon)$

Theorem  $E[R_T] \leq \frac{\log k}{\varepsilon} + \varepsilon \sum_{t=1}^T p_t^T l_t^2$

$\underbrace{\sum_{x=1}^k p_t(x) l_t(x)^2}_{\leq \varepsilon \cdot T \text{ if } l_t(x) \in [0, 1]}$

for  $\varepsilon := \frac{\sqrt{\log k}}{\sqrt{T}}$  Assuming  $l_t \in [0, 1]$   $\Rightarrow E[R_T] \leq 2\sqrt{T \log k}$

$\Rightarrow E[\frac{R_T}{T}] \leq \frac{2\sqrt{\log k}}{\sqrt{T}}$

Proof  $\phi_t = \sum_{x=1}^k w_t(x)$ . Then  $\phi_1 = k$

$$\begin{aligned}\phi_{t+1} &= \sum_{x=1}^k w_t(x) \cdot \exp(-\varepsilon l_t(x)) \\ &= \phi_t \sum_{x=1}^k p_t(x) \exp(-\varepsilon l_t(x)) \\ &\leq \phi_t \sum_x p_t(x) (1 - \varepsilon l_t(x) + \varepsilon^2 l_t^2(x)) \\ &= \phi_t (1 - \varepsilon p_t^\top l_t + \varepsilon^2 p_t^\top l_t^2) \\ &\leq \phi_t \exp(-\varepsilon p_t^\top l_t + \varepsilon^2 p_t^\top l_t^2) \\ &\leq \underbrace{\phi_1}_{k} \exp(-\varepsilon \sum_{t=1}^T p_t^\top l_t + \varepsilon^2 \sum_{t=1}^T p_t^\top l_t^2)\end{aligned}$$

Note:  $p_t(x) = \frac{w_t(x)}{\phi_t}$

$e^{-x} \leq 1 - x + x^2$

for  $x \geq 0$

$e^x \geq 1 + x \quad \forall x$

For best expert  $x^*$ ,  $w_T(x^*) = \exp(-\varepsilon \sum_{t=1}^T l_t(x^*))$

Thus:  $w_T(x^*) \leq \phi_T \leq k \cdot \exp(-\varepsilon \sum_{t=1}^T p_t^\top l_t + \varepsilon^2 \sum_{t=1}^T p_t^\top l_t^2)$

$\log$   $\Rightarrow -\varepsilon \sum_{t=1}^T l_t(x^*) \leq \log k - \varepsilon \sum_{t=1}^T p_t^\top l_t + \varepsilon^2 \sum_{t=1}^T p_t^\top l_t^2$

$$\Rightarrow \mathbb{E}[R_T] = \sum_{t=1}^T p_t^\top l_t - \sum_{t=1}^T l_t(x^*) \leq \frac{\log k}{\varepsilon} + \varepsilon \sum_{t=1}^T p_t^\top l_t^2$$



## BANDIT SETTING

DM only observes  $l_t(x_t)$   
 $\hookrightarrow$  exploration-exploitation

key idea: Reduction to full information  
 Setting by constructing  
 estimates of  $l_t(x)$

Ex. Recommender systems

$X = \text{set of } k \text{ possible recommend.}$   
 At time  $t$ , user arrives  
 pick  $x_t$ , if user interested  
 in  $x_t$ , user watches/does/-.  
 $l_t(x_t) = \begin{cases} 0 & \text{if user interested} \\ 1 & \text{oth.} \end{cases}$

Sps we play  $x_t \sim p_t$  (as in Hedge)

Define  $\tilde{l}_t(x) := \begin{cases} \frac{1}{p_t(x_t)} l_t(x_t) & \text{if } x = x_t \\ 0 & \text{oth.} \end{cases}$

$$\text{Then } \forall x: \mathbb{E}_{x_t \sim p_t} [\tilde{l}_t(x)] = \sum_{x_t} p_t(x_t) \cdot \tilde{l}_t(x_t) = p_t(x_t) \cdot \tilde{l}_t(x_t) + 0 \\ = p_t(x_t) \cdot \frac{l_t(x_t)}{p_t(x_t)} = l_t(x_t)$$

$$\forall x: \mathbb{E}_{x_t \sim p_t} [\tilde{l}_t^2(x)] = p_t(x_t) \cdot \tilde{l}_t^2(x_t) = p_t(x_t) \frac{l_t^2(x_t)}{p_t^2(x_t)} = \frac{l_t^2(x_t)}{p_t(x_t)}$$

Algorithm: EXP3 : Play Hedge on  $\tilde{l}_t$

Theorem: For  $l_t \in [0, 1]$   $\forall t, k$ : Sps we ran EXP3 w.  $E = \sqrt{\frac{\log k}{T}}$

Then it holds:  $\mathbb{E}[R_T] \leq 2 \sqrt{T \log k}$  due to partial inf.

Proof:

$$\mathbb{E}[R_T] = \mathbb{E}\left[\sum_{t=1}^T l_t(x_t) - \sum_{t=1}^T l_t(x^*)\right]$$

$$= \mathbb{E}\left[\sum_{t=1}^T p_t^\top l_t - \sum_{t=1}^T l_t(x^*)\right]$$

$$= \mathbb{E}\left[\sum_{t=1}^T p_t^\top \tilde{l}_t - \sum_{t=1}^T \tilde{l}_t(x^*)\right] \quad \begin{array}{l} \text{lin. of expect.} \\ x^* \text{ independent of} \\ \tilde{l}_t; \quad \mathbb{E}[\tilde{l}_t] = l_t \end{array}$$

Hedge on

$$\tilde{l}_{1:T} \subseteq \mathbb{E}\left[\epsilon \cdot \sum_{t=1}^T p_t^\top \tilde{l}_t^2 + \frac{\log k}{\epsilon}\right]$$

Note:  $\mathbb{E}[p_t^\top \tilde{l}_t^2] = \sum_{k=1}^k p_t(k) \cdot \mathbb{E}_{x_t \sim p_t} [\tilde{l}_t^2(k)]$

$$= \sum_{k=1}^k p_t(k) \frac{\tilde{l}_t^2(k)}{p_t(k)} = \sum_{k=1}^k \tilde{l}_t^2(k) \leq k$$

$$\Rightarrow \mathbb{E}[R_T] \leq \epsilon \cdot k \cdot T + \frac{\log k}{\epsilon}$$

$$\stackrel{\text{Def. of } \epsilon}{\leq} 2 \sqrt{T \cdot k \cdot \log k}$$

## LEARNING IN REPEATED GAMES

For  $t = 1: T$

We pick  $x_t \in X$

Opp. picks  $y_t \in Y$

$$\text{We obtain } l_t^X(x_t) = f(x_t, y_t)$$

$$\text{Opp. obtains } l_t^Y(y_t) = 1 - f(x_t, y_t)$$

fixed

Can directly apply Hedge (or EXP3) with sublinear regret  $O(\sqrt{T \log k})$  (or  $O(\sqrt{T k \log k})$ )

Fact if both players play no-regret, the "average actions"  $P_T^X = \text{Unif}\{x_1 \dots x_T\}$   
 $P_T^Y = \text{Unif}\{y_1 \dots y_T\}$

Converge to a Nash equilibrium in game f

## Outlook

Exponential gap (in  $\delta$ ) between full info and bandit setting.

In hindsight, often can also observe  $y_t$  (other player's action)

Can show:

$$\mathbb{E}[R_T] = O\left(\sqrt{T \log k} + \gamma_T \sqrt{T}\right)$$

$P$   
regret under  
full info

$\uparrow$   
depends on  
structure of game

E.g.  $f(x, y) = w^\top \phi(x, y)$ , for  $w \in \mathbb{R}^d$

$$\bar{l}_t(x) \Rightarrow \gamma_T = d \cdot \log T$$

