

# Lecture 09

Other ERM algs:

[skip MWA is ERM]

e.g. Sampled fictitious play, Li, Tewari, 2017  
(is Hannon consistent)

FTPL, etc.

Last time:

$$(ERM) \quad \limsup_{T \rightarrow \infty} \frac{1}{T} \max_{a_i \in A_i} \left( \sum_{t=0}^{T-1} U_i(a_i, a_{-i}^t) - U_i(a_i^t, a_{-i}^t) \right) \leq 0 \quad \text{a.s.}$$

ERM<sub>i</sub>

randomization by all players

• If all players use ERM strategies, then a.s.  

$$\limsup_{T \rightarrow \infty} \sum_{a \in A} q^T(a) (U_i(a_i, a_{-i}) - U(a)) \leq 0 \quad (\text{limiting behaviors})$$

⇒ Coarse Correlated Equilibrium:

$$CCE := \left\{ q \in \Delta(A) : \sum_{a \in A} q(a) U_i(a_i, a) \leq \sum_{a \in A} q(a) U_i(a), \forall i, a_i \right\}$$

• Recall

$$CE := \left\{ q \in \Delta(A) : \sum_{a_{-i} \in A_{-i}} q(a_{-i} | a_i) U_i(a_i, a_{-i}) \leq \sum_{a_{-i} \in A_{-i}} q(a_{-i} | a_i) U_i(a_{-i}), \forall i, a_i, a_i' \right\}$$

$$\Leftrightarrow CE := \left\{ q \in \Delta(A) : \sum_{a_{-i} \in A_{-i}} q(a_{-i}, a_i) U_i(a_i', a_{-i}) \leq \sum_{a_{-i} \in A_{-i}} q(a_{-i}, a_i) U_i(a_i), \forall i, a_i, a_i' \right\}$$

•  $CE \subseteq CCE$   
 ? ERM

Consider the following Internal Regret Minimization:

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \max_{a_i, a_i' \in A_i} \sum_{t=0}^{T-1} \mathbb{I}(a_i^t = a_i) (U_i(a_i', a_{-i}^t) - U_i(a_i, a_{-i}^t)) \text{ a.s.}$$

(replace every  $a_i^t = a_i$  by  $a_i'$ )

• If all players use IRM strategies, then

$$\limsup_{T \rightarrow \infty} \sum_{a_i \in A_i} q(a_i) (U_i(a_i, a_{-i}) - U_i(a_{-i}, a_i)) \leq 0$$

⇒ Theorem 9.1

If all players use IRM strategies, then the joint empirical distribution approaches the set of correlated equilibria (a.s.).

• What strategies are IRM?

• Consider this generalized/unified regret definition:

$$R_{\mathcal{F}} := \max_{F \in \mathcal{F}} \{ U_{\pi, F}^T - U_{\pi} \}$$

Modification rule

Given a mixed strategy  $p^+$  from  $\pi$ , and a modification rule  $F^t$ :

$$f^t = F^t(p^+)$$

$$\text{where } f_i^t = \sum_{j: F^t(j)=i} p_j^t$$

$A \rightarrow A$

$$\cdot F := \{ F_t : t \geq 1 \}$$

•  $U_{\pi, F}^T$  is the profit induced by  $(f^t : t \geq 1)$ .

≡ Recall external regret:

$$ER := \underbrace{U_{\theta, \max}^T}_{\max_{\theta \in \Theta} U_{\theta}^T} - U_{\pi}^T$$

$$\max_{\theta \in \Theta} U_{\theta}^T$$

comparison class of algs

equivalent to  $F^{ex} := \{ F_i : i \in A \}$  where  $F_i(j) = i \forall i, j \in A$

≡ Internal regret =

$$IR := \max_{i, j \in A} \left\{ \sum_{t=1}^T p_i^t (u_j^t - u_i^t) \right\}$$

fixing other players' strategies  
mixed strategies by  $\pi$

Q? equivalent to  $F^{in} := \left\{ F_{ij} : \begin{cases} F_{ij}(i) = j \\ F_{ij}(i') = i' \forall i' \neq i \end{cases} \cdot \begin{matrix} i, j \in A \\ i \neq j \end{matrix} \right\}$   
 $|F^{in}| = |A|(|A|-1)$

- In general, how many distinct  $F: A \rightarrow A$  are there?

$$|A|^{|A|}$$

- Each  $F$  swaps the current online action  $i$  w/  $F(i)$

$$\begin{aligned} \text{Swap regret: } SR &:= \max_{F \in \mathcal{F}^{\text{sw}}} \left\{ U_{\pi, F}^T - U_{\pi} \right\} \\ &= \sum_{i=1}^{|A|} \max_{j \in A} \left\{ \sum_{t=1}^T p_i^t (u_j^t - u_i^t) \right\} \end{aligned}$$

- We see examples of ERM policies, such as MWA, sampled FP, leading to CCE.

What are IRM policies for CE?

- A general reduction from ER to IR (in fact, SR).

[Blum & Mansour, COLT 2005]

Create  $N = |A|$  copies of the external-regret procedures

$\pi_1, \dots, \pi_N$ , each of them satisfying

not exactly ERM, can generalize

$$\forall \text{ sequence of } T \text{ payoffs (losses)} (u^t: t=1 \dots T), \forall j \in \{1, \dots, N\}.$$

$$U_{\pi_i}^T = \sum_{t=1}^T u_{\pi_i}^t \leq \sum_{t=1}^T u_j^t + R = U_j^T + R$$

We combine the  $N$  procedures to one master procedure as

follows. At each time  $t$ , each procedure  $\pi_i$  outputs a distribution

$$q_i^t = (q_{ij}^t: j=1, \dots, N).$$

fraction it assigns to action  $j$   $\{q_{ij}^t: i, j \in A\}$

$$\text{Compute } p_j^t = \sum_i p_i^t q_{ij}^t \iff p^t = Q^t p^t$$

a stationary distribution of the Markov Process defined by  $Q^t$ .

[ exercise: show this  $p^+$  always exists ]

Choosing  $p^+$  can be regarded in two equivalent ways:

- (1). Using  $p^+$  to select action  $j$  w/ probability  $p_j^+$  (actual payoff)  $p^+ \cdot u^+$
- (2). Select procedure  $\pi_i$  w/ probability  $p_i^+$  and then use  $\pi_i$  to select the action  $(Q^+ p^+)$

After receiving a payoff (full information model)  $u^+$ , we return to each  $\pi_i$  the vector  $p_i^+ u^+$ . So, procedure  $\pi_i$  experiences

$$(p_i^+ u^+) \cdot q_i^+ = p_i^+ (q_i^+ \cdot u^+)$$

inner product

Since  $\pi_i$  is an  $R$  external regret procedure,  $\forall$  action  $j$ ,

$$\sum_{t=1}^T p_i^t (q_i^t \cdot u^t) \leq \sum_{t=1}^T p_i^t u_j^t + R \quad (*)$$

Summing the payoffs of the  $N$  procedures, at a given time  $t$ ,

$$\sum_i p_i^t (q_i^t \cdot u^t) = \underbrace{(Q^+ p^+)}_{p^+} \cdot u^t = \underbrace{p^+ \cdot u^t}_{\text{actual payoff}}$$

Hence, (\*) gives (summing over  $i=1, \dots, N$ )

$$u_{\pi}^T = \sum_{t=1}^T p^+ \cdot u^t$$

$$u_{\pi}^T \leq \sum_{i=1}^N \sum_{t=1}^T p_i^t u_{F(i)}^t + NR = u_{\pi, F}^T + NR$$

for any function  $F: \underbrace{\{1, \dots, N\}}_A \rightarrow \underbrace{\{1, \dots, N\}}_A$ .

In summary.

### Theorem 9.2

Given an  $R$  external regret procedure, the constructed master online procedure  $\pi$  has the following guarantee.  $\forall F: A \rightarrow A$ ,

$$U_{\pi} \leq U_{\pi, F} + NR$$

i.e., the swap regret of  $\pi$  is at most  $NR$ .

Hence.

### Corollary 9.1

There exists an online algorithm  $\pi$  s.t. for every function  $F: A \rightarrow A$ , we have that

combining MWA. for example

$$U_{\pi} \leq U_{\pi, F} + O(N\sqrt{T \log N})$$

i.e., the swap regret of  $\pi$  is at most  $O(N\sqrt{T \log N})$ .

- Read Section 4.5, 4.6 in the book, which can be found in our reading materials for more detailed discussions on the partial information setting.

• More on online learning algs.

• Recall: MWA assumptions:

(1) The set of allowed actions for the player is the probability simplex

$$K_N = \left\{ P \in \mathbb{R}_+^N, \sum_{i=1}^N P_i = 1 \right\}$$

(2) Loss / Payoff functions are linear,  $f^i(P) = m^i \cdot P$  where  $m^i \in \mathbb{R}^N$  is normalized s.t.  $|m_i^j| \leq 1 \forall i \in \{1, \dots, N\}$ . This ensures  $f^i(P) \in [-1, 1], \forall P$ .

Given this, MWA is a special instance of

Follow the Regularized leader (FTRL)

$$P^{t+1} = \arg \min_{P \in K_N} \left\{ \eta \sum_{j \text{ st}} m^j \cdot P + R(P) \right\} \quad \begin{array}{l} -H(P) \\ \text{"} \end{array}$$

• Now, let's turn to online convex optimization (OCO):

Goal: solve  $\min_{x \in K} \sum_{t=1}^T f_t(x)$  online

•  $K$  is bounded, convex, closed

•  $f_t: K \rightarrow \mathbb{R}$  is convex.

Similar to our previous discussion for MWA, the FTL / fictitious play scheme below fails in the worst-case:

$$x_{t+1} = \arg \min_{x \in K} \sum_{\tau=1}^t f_{\tau}(x)$$

• Consider  $K = [-1, 1]$ ,  $f_1(x) = \frac{1}{2}x$ ,  $f_{\tau}(x)$  for  $\tau=2, \dots, T$  alternate between  $-x$  and  $x$ . Thus,

$$\sum_{\tau=1}^t f_{\tau}(x) = \begin{cases} \frac{1}{2}x & t \text{ is odd} \\ -\frac{1}{2}x & t \text{ is even} \end{cases}$$

$\Rightarrow$  FTL strategy keeps shifting between  $x_t = -1$  and  $x_t = 1$ , making the wrong choices.

- Consider regularization functions  $R: K \rightarrow \mathbb{R}$  that are strongly convex, smooth, and twice differentiable.

Hence, by strong convexity, the Hessian  $\nabla^2 R(x) \succ 0$  is positive definite.

Define the diameter of  $K$  as

$$D_R := \left( \max_{x, y \in K} \{ R(x) - R(y) \} \right)$$

- Dual norm of  $\|\cdot\|$ :  $\|y\|^* := \sup_{\|x\| \leq 1} \{x^T y\}$ .

Dual norm of the matrix norm  $\|x\|_A = \sqrt{x^T A x}$ :  $\|x\|_A^* = \|x\|_{A^{-1}}$

For notational simplicity, we write Generalized Cauchy-Schwarz inequality:  $x^T y \leq \|x\|_A \cdot \|y\|_A^*$   
↓  
any norm

$$\|x\|_y := \|x\|_{\nabla^2 R(y)}$$

$$\|x\|_y^* := \|x\|_{\nabla^{-2} R(y)}$$

- Difference between the value of the regularization function at  $x$  and the value of the 1st order Taylor approximation:

$$\text{Bregman divergence: } B_R(x \| y) := R(x) - R(y) - \nabla R(y)^T (x - y).$$

For twice differentiable functions, Taylor expansion and the mean-value theorem implies

$$B_R(x \| y) = \frac{1}{2} \|x - y\|_{\mathbb{Z}}^2 \text{ for some } \mathbb{Z}, \text{ and } \exists \alpha \in [0, 1] \text{ s.t.}$$

$$\mathbb{Z} = \alpha x + (1 - \alpha)y.$$

Thus, the Bregman divergence defines a local norm, which has a dual norm  $\|\cdot\|_{x, y}^* := \|\cdot\|_{\mathbb{Z}}^*$

Finally, we write  $\|\cdot\|_+ = \|\cdot\|_{x_t, x_{t+1}}$  so that  $B_R(x_t \| x_{t+1}) = \frac{1}{2} \|x_t - x_{t+1}\|_+^2$

To wrap up, let's summarize and revisit this FTRL meta-algorithm:

### Algorithm FTRL

Input:  $\eta > 0$ , regularization function  $R$ , and  $K$

$$\text{let } x_1 = \underset{x \in K}{\operatorname{argmin}} \{ R(x) \}$$

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

play  $x_t$  and receive cost  $f_t(x_t)$  (can relax from the full information setting to the bandit gradient setting,  $\nabla_t := \nabla_{f_t}(x_t)$ )

update

$$x_{t+1} = \underset{x \in K}{\operatorname{argmin}} \left\{ \eta \sum_{s=1}^t f_s(x) + R(x) \right\}$$

relax to  $\nabla_s^T x$   
since  $\forall t$

end for

$$(\alpha) \leftarrow f_t(x_t) - f_t(x^*) \leq \nabla_t^T (x_t - x^*)$$

• Theoretical Guarantee:

by convexity

### Theorem 9.3

The FTRL algorithm attains  $\forall u \in K$  the following regret bound:

$$ER(\text{FTRL}) \leq 2\eta \sum_{t=1}^T (\|\nabla_t\|_+^*)^2 + \frac{R(u) - R(x_1)}{\eta}$$

• Note that if  $\|\nabla_t\|_+^* \leq G_R \forall t$ , then optimizing over  $\eta$  gives a bound of  $2D_R G_R \sqrt{2T}$ .

Proof: Lemma: FTRL guarantees  $ER(\text{FTRL}) \leq \sum_{t=1}^T \nabla_t^T (x_t - x_{t+1}) + \frac{1}{\eta} P_R^2$ .

↓  
proof: Define  $g_0(x) := \frac{1}{\eta} R(x)$ .

$$g_t(x) := \nabla_t^T x$$

By (α), it suffices to bound  $\sum_{t=1}^T [g_t(x) - g_t(u)]$ . We

first show that  $\forall u \in K$ ,

$$\sum_{t=0}^T g_t(u) \geq \sum_{t=0}^T g_t(x_{t+1})$$

To see it, we use induction on  $T$ :



• Induction base: by definition,  $x_1 = \operatorname{argmin}_{x \in K} R(x)$ , thus  $g_0(u) \geq g_0(x_1) \forall u$ .

• Induction step:

$$\text{Assume for } T. \text{ we have } \sum_{t=0}^T g_t(u) \geq \sum_{t=0}^T g_t(x_{t+1})$$

$$\text{For } T+1, \text{ since } x_{T+2} = \operatorname{argmin}_{x \in K} \left\{ \sum_{t=0}^{T+1} g_t(x) \right\},$$

$$\begin{aligned} \sum_{t=0}^{T+1} g_t(u) &\geq \sum_{t=0}^{T+1} g_t(x_{T+2}) \\ &= \sum_{t=0}^T g_t(x_{T+2}) + g_{T+1}(x_{T+2}) \\ &\geq \sum_{t=0}^T g_t(x_{t+1}) + g_{T+1}(x_{T+2}) \\ &= \sum_{t=0}^{T+1} g_t(x_{t+1}) \end{aligned}$$

induction hypothesis  
 $u = x_{T+2}$

As a conclusion,

$$\begin{aligned} \sum_{t=1}^T [g_t(x_t) - g_t(u)] &\leq \sum_{t=1}^T [g_t(x_t) - g_t(x_{t+1})] + [g_0(u) - g_0(x_1)] \\ &= \sum_{t=1}^T [g_t(x_t) - g_t(x_{t+1})] + \frac{1}{\eta} [R(u) - R(x_1)] \\ &\leq \sum_{t=1}^T [g_t(x_t) - g_t(x_{t+1})] + \frac{1}{\eta} D_R^2 \quad \# \end{aligned}$$

Now, since  $R(x)$  is a convex function and  $K$  is a convex set.

$$\text{Define } \Phi_t(x) := \left( \sum_{s=1}^t \nabla_s^T x + R(x) \right).$$

The Taylor expansion implies

$$\begin{aligned} \Phi_t(x) &= \Phi_t(x_{t+1}) + (x - x_{t+1})^T \nabla \Phi_t(x_{t+1}) + B_{\Phi_t}(x \parallel x_{t+1}) \\ &\geq \Phi_t(x_{t+1}) + B_{\Phi_t}(x \parallel x_{t+1}) \quad \text{since } x_{t+1} \text{ minimizes } \Phi_t \text{ over } K \\ &= \Phi_t(x_{t+1}) + B_R(x \parallel x_{t+1}) \quad \text{since the term } \nabla_s^T \cdot x \text{ is linear, it won't affect the Bregman divergence.} \end{aligned}$$

Rearrange the terms. We get

$$\begin{aligned} B_R(x_t \| x_{t+1}) &\leq \bar{\Phi}_+(x_t) - \bar{\Phi}_+(x_{t+1}) \\ &= (\bar{\Phi}_{t+1}(x_t) - \bar{\Phi}_{t+1}(x_{t+1})) + \eta \nabla_+^T(x_t - x_{t+1}) \\ &\leq \eta \nabla_+^T(x_t - x_{t+1}) \quad (x_t \text{ minimizes } \bar{\Phi}_{t+1} \text{ again}) \end{aligned}$$

Using our notation, the generalized Cauchy-Swartz inequality.

$$\begin{aligned} \nabla_+^T(x_t - x_{t+1}) &\leq \|\nabla_+\|_+^* \cdot \|x_t - x_{t+1}\|_+ \\ &= \|\nabla_+\|_+^* \cdot (2B_R(x_t \| x_{t+1}))^{\frac{1}{2}} \\ &\leq \|\nabla_+\|_+^* \cdot (2\eta \nabla_+^T(x_t - x_{t+1}))^{\frac{1}{2}} \end{aligned}$$

$$\Rightarrow \nabla_+^T(x_t - x_{t+1}) \leq 2\eta (\|\nabla_+\|_+^*)^2$$

Substituting this into the lemma complete the proof of Theorem 9.3. #

• connection to online mirror descent (OMD)

OMD is a general class of 1st order methods extending GD.

It has two versions, agile and lazy.

### Algorithm OMD

Input:  $\eta > 0$ , regularization function  $R(x)$ .

Let  $y_1$  be s.t.  $\nabla R(y_1) = 0$

Let  $x_1 = \operatorname{argmin}_{x \in K} B_R(x \| y_1)$

for  $t = 1, \dots, T$ . do

play  $x_t$

observe the loss  $f_t$ , let  $\nabla_+ = \nabla f_t(x_t)$

update  $y_t$  according to

$$\nabla R(y_{t+1}) = \nabla R(y_t) - \eta \nabla_+ \quad \text{Lazy}$$

$$\nabla R(y_{t+1}) = \nabla R(x_t) - \eta \nabla_+ \quad \text{Agile}$$

Project according to  $B_R$ :

$$x_{t+1} = \operatorname{argmin}_{x \in K} B_R(x \| y_{t+1})$$

end for

- Both two version have regret bound guarantees similar to FTRL
- For instance,

Theorem 9.4 (equivalence between FTRL and lazy OMD)

Let  $f_1, \dots, f_T$  be linear cost functions. The lazy OMD and FTRL are identical

i.e.,

$$x_t = \operatorname{argmin}_{x \in K} B_R(x \| y_t) = \operatorname{argmin}_{x \in K} \left( \eta \sum_{s=1}^{t-1} \nabla_s^T \cdot x + R(x) \right)$$

Proof: The optimal solution

$x_t^*$  of the unconstrained optimization (0) satisfies

$$\nabla R(x_t^*) = -\eta \sum_{s=1}^{t-1} \nabla_s$$

By the lazy OMD update rule:

$$\nabla R(y_t) = -\eta \sum_{s=1}^{t-1} \nabla_s$$

$$\Rightarrow \nabla R(x_t^*) = \nabla R(y_t)$$

Since  $R$  is strictly convex,  $x_t^* = y_t$

$$\begin{aligned} \text{Hence, } B_R(x \| y_t) &= R(x) - R(y_t) - (\nabla R(y_t))^T (x - y_t) \\ &= R(x) - R(y_t) + \eta \sum_{s=1}^{t-1} \nabla_s^T \cdot (x - y_t) \end{aligned}$$

$\Rightarrow$  It's equivalent to

$$\text{minimize } R(x) + \eta \sum_{s=1}^{t-1} \nabla_s^T \cdot x \text{ over } K \quad \#$$

independent of  $x$

- Now, we apply the general regret bound in Theorem 9.3 to concrete examples of  $R(x)$ .

Case I:  $R(x) = x \log x$

$$\Rightarrow \nabla R(x) = 1 + \log x$$

$$K = \left\{ x \in \mathbb{R}_+^n : \sum x_i = 1 \right\} \Rightarrow x_{t+1}(i) = \frac{x_t(i) \cdot e^{-\eta \nabla_t(i)}}{\sum_{j=1}^n x_t(j) \cdot e^{-\eta \nabla_t(j)}}$$

If costs are in  $[-1, 1]$ .

$$\|\nabla_t\|_t^* \leq \|\nabla_t\|_\infty \leq 1 =: G_R$$

The diameter satisfies  $D_R^2 \leq \log \eta$

$$\Rightarrow ER(\text{MWA}) \leq 2D_R G_R \sqrt{2T} \leq 2\sqrt{2T \log \eta}.$$

Case II  $R(x) = \frac{1}{2} \|x - x_0\|_2^2$ .  
arbitrary  $x_0 \in K$

$$\nabla R(x) = x - x_0$$

$$\Rightarrow x_t = \text{Proj}_K(y_t), \quad y_t = y_{t-1} - \eta \nabla_{t-1} \quad (\text{lazy})$$

$$x_t = \text{Proj}_K(y_t), \quad y_t = x_{t-1} - \eta \nabla_{t-1} \quad (\text{agile})$$

exactly online GD.

$$\begin{aligned} \Rightarrow ER(\text{OGD}) &\leq \frac{1}{\eta} D_R^2 + 2\eta \sum_{t=1}^T (\|\nabla_t\|_t^*)^2 \\ &\leq \frac{1}{2\eta} D^2 + 2\eta \sum_{t=1}^T \|\nabla_t\|^2 \quad (\text{HW } \|\cdot\|_t \text{ reduces to } \|\cdot\|_2) \end{aligned}$$

$$\leq 2GD\sqrt{T}$$

$\max_{x, y \in K} \|x - y\|$  Euclidean diameter

$$\|\nabla f_t(x)\| \leq G, \quad \forall t, x \in K$$

• What about nonconvex costs?

Read the book by Elad Hazan for more detailed discussions.