

# Lecture 10: Correlated Equilibria and No Regret Learning

Jun Wang  
UCL

# Table of content

- ▶ Correlated Equilibria and Coarse Correlated Equilibria
  - ▶ Definitions
- ▶ Online Learning
  - ▶ Setting
  - ▶ Performance Evaluation
  - ▶ Regret and No-regret Algorithms
- ▶ Connection to games
  - ▶ Why online learning and no-regret algorithms?
  - ▶ Connection to Equilibria
  - ▶ A No-Regret Algorithm: Multiplicative Weights
  - ▶ Correlated Equilibria and Swap Regret

# Main references

- ▶ Chapters 17-18, Roughgarden, Tim. Twenty lectures on algorithmic game theory. Cambridge University Press, 2016. [Rou16]
- ▶ Chapter 2.2, Fudenberg, Drew, and Jean Tirole. "Game theory" Cambridge, Massachusetts 393.12 (1991): 80. [FT91]
- ▶ Blum, Avrim, and Yishay Mansour. "Learning, regret minimization, and equilibria." (2007). [BM07]
- ▶ Cesa-Bianchi, Nicolo, and Gábor Lugosi. Prediction, learning, and games. Cambridge university press, 2006. [CL06]
- ▶ Lecture notes from 601.436/636 Algorithmic Game Theory - Spring 2020 by Michael Dinitz
- ▶ Asu Ozdaglar, 6.254 : Game Theory with Engineering Applications Lecture 4: Strategic Form Games - Solution Concepts

## A motivation example

- ▶ Nash equilibrium may act as a minimal necessary condition for *reasonable* predictions in the situation where the players must choose their strategies *independently*.
- ▶ Alternatively, players may engage in *preplay* discussion, and then go off to isolated rooms to choose their strategies.
  - ▶ In some cases, both players might gain if they build a "signalling device" that sent signals to the separate rooms<sup>1</sup>.

	Ballet	Football
Ballet	1,4	0,0
Football	0,0	4,1

- ▶ Suppose that the players flip a coin and go to the Ballet if the coin is Heads, and to the Football if the coin is tails,
  - ▶ i.e., they randomize between two pure strategy Nash equilibria, resulting in a payoff of  $(5/2, 5/2)$  that is not a Nash equilibrium payoff.

---

<sup>1</sup>Robert J Aumann. "Subjectivity and correlation in randomized strategies".

# Correlated Equilibria

Consider a cost-minimisation game<sup>2</sup> with  $k$  players  $[k]$  and strategy sets  $\{S^i\}_{i \in [k]}$ , with  $S = S^1 \times S^2 \times \dots \times S^k$  and cost functions  $c^i : S \rightarrow \mathbb{R}$  for each player  $i \in [k]$ .

Let us think a new way of playing this game:

- ▶ There is a “trusted third party”  $U$  recommending all players playing according to a publicly known  $\sigma$  which is a distribution over  $S = S^1 \times S^2 \times \dots \times S^k$ .
- ▶  $U$  samples some strategy profile  $s \in S$  from  $\sigma$  (but keeps  $s$  secret).
- ▶ For each player  $i \in [k]$ ,  $U$  privately tells  $s^i = s[i]$  to  $i$ .
- ▶ Now each player  $i$  decides whether to play  $s^i$  or to deviate to some other strategy.

---

<sup>2</sup>The formalisms of cost-min and payoff-max games are equivalent; we use cost-min to be consistent with online learning.

# Correlated Equilibria

## Definition

Let  $\sigma$  be a distribution over  $S = S^1 \times \dots \times S^k$ . Then  $\sigma$  is a *correlated equilibrium* if

$$\mathbb{E}_{s \sim \sigma} [c^i(s) \mid s^i] \leq \mathbb{E}_{s \sim \sigma} [c^i(s^{-i}, s'') \mid s^i],$$

for all  $i \in [k]$  and for all  $s'' \in S^i$ .

- ▶ Each player  $i$  knows  $\sigma$  and  $s^i$ , so can figure out their expected cost if no one deviates:  $\mathbb{E}_{s \sim \sigma} [c^i(s) \mid s^i]$ .
- ▶ Player  $i$  can also figure out their expected cost if they deviated to  $s''$  but no one else deviates:  $\mathbb{E}_{s \sim \sigma} [c^i(s^{-i}, s'') \mid s^i]$ .
- ▶ So the correlated equilibrium condition says that no one has incentive to deviate when the game is played this way.
- ▶ every Nash equilibrium is a correlated equilibrium, since if  $\sigma$  is a product distribution then the conditioning affects nothing.
- ▶ But there can be correlated equilibria which are not Nash.

## Coarse Correlated Equilibria<sup>3</sup>

- ▶ CCE: the mediator requires more commitment from the players.
- ▶ It asks the players, before *running the lottery* (a joint distribution of actions  $\sigma$ ), to either commit to the future outcome of the lottery or play any strategy of their own without learning anything about the outcome of the lottery.
- ▶ The equilibrium property is that each player finds it optimal to commit before it happens to use the strategy selected by the lottery.

---

<sup>3</sup>Hervé Moulin and J-P Vial. "Strategically zero-sum games: the class of games whose completely mixed equilibria cannot be improved upon". In: *International Journal of Game Theory* 7.3-4 (1978), pp. 201–221.

# Coarse Correlated Equilibria

## Definition

Let  $\sigma$  be a distribution over  $S = S^1 \times \dots \times S^k$ . Then  $\sigma$  is a *coarse correlated equilibrium* if

$$\mathbb{E}_{s \sim \sigma} [c^i(s)] \leq \mathbb{E}_{s \sim \sigma} [c^i(s^{-i}, s'')],$$

for all  $i \in [k]$  and for all  $s'' \in S^i$ .

- ▶ Like correlated equilibria we are now allowed to have non-product distributions  $\sigma$ ,
- ▶ but like Nash equilibria there is no conditioning on our action.
- ▶ It is the same setup as correlated equilibria with their being a trusted third party, but now each player  $i$  has to decide whether to deviate before being told  $s^i$ .





# Coarse Correlated Equilibria

## Definition

Let  $\sigma$  be a distribution over  $S = S^1 \times \cdots \times S^k$ . Then  $\sigma$  is a coarse correlated equilibrium if

$$\mathbb{E}_{s \sim \sigma} [c^i(s)] \leq \mathbb{E}_{s \sim \sigma} [c^i(s^{-i}, s'')],$$

for all  $i \in [k]$  and for all  $s'' \in S^i$ .

- ▶ It's not hard to prove (good exercise to do at home!) that any correlated equilibrium is a coarse correlated equilibrium.
- ▶ The converse is definitely not true, though: there are coarse correlated equilibria that are not correlated equilibria.
  - ▶ Even if we don't have incentive to deviate if we're not told what to play, we might have incentive to deviate if there are particular actions that we're told to play.
- ▶  $CE \subset CCE$

How do we learn CE and CCE?

# Online Learning

## Setting:

- ▶ A player  $i$  who has an action set  $\mathbb{A}$  where  $n = |\mathbb{A}|$  (also known as a set of arms in *multi-armed bandit* problems related literature).
- ▶ At time  $t = 1, 2, \dots, T$ :
  - ▶ The player  $i$  picks a mixed strategy  $\pi_t \in \Delta(\mathbb{A})$  where  $\Delta$  is a probability simplex.
  - ▶ The player's adversary picks a cost vector  $c_t : \mathbb{A} \rightarrow [0, 1]$ .
  - ▶ An action  $a_t$  is drawn from  $\pi_t$ , and the algorithm has cost  $c_t(a_t)$ . The player learns either
    - ▶ the entire vector  $c_t$  (in the “experts” setting) or
    - ▶ just learns  $c_t(a_t)$  (in the “bandit” setting).

# How do we evaluate a player's performance online?

The best action sequence in hindsight:

$$OPT = \sum_{t=1}^T \min_{a \in \mathbb{A}} c_t(a)?$$

## Theorem

*No algorithm can be competitive with the best action sequence in hindsight.*

## Proof.

Let  $|\mathbb{A}| = 2$ . The adversary will work as follows: if  $\pi_t(0) \geq \frac{1}{2}$  then the adversary sets  $c_t(0) = 1$  and  $c_t(1) = 0$  and otherwise  $\pi_t(1) \geq \frac{1}{2}$  and the adversary sets  $c_t(1) = 1$  and  $c_t(0) = 0$ . Then there is some sequence that has total cost 0, but at every time the algorithm has probability at least  $\frac{1}{2}$  of getting cost 1. Thus the expected cost of the algorithm is at least  $\frac{T}{2}$  while  $OPT = 0$ , for an unbounded competitive ratio.  $\square$

# How do we evaluate a player's performance online?

We can, however, compete with the *best single action in hindsight*:

$$\min_{a \in \mathbb{A}} \sum_{t=1}^T c_t(a)$$

## Definition

The (external) regret of an action sequence  $\{a_1, a_2, \dots, a_T\}$  with respect to an arbitrary action  $a \in \mathbb{A}$  is

$$R_T(a) = \frac{1}{T} \left( \sum_{t=1}^T c_t(a_t) - \sum_{t=1}^T c_t(a) \right)$$

- ▶ One possible motivation: we consider each of these strategies to be an expert that gives us advice.
- ▶ The mentioned benchmark is cost that the best of these experts incurs, so if we had relied on only this expert's advice.

# What do we know about our adversary?

## Definition

An (adaptive) adversary is a function which take as input 1) an algorithm  $A$ , 2) the time  $t$ , 3) mixed strategies  $\{\pi_1, \pi_2, \dots, \pi_t\}$  produced by  $A$  and 4) realized actions  $\{a_1, a_2, \dots, a_{t-1}\}$  from the past, and outputs a cost vector  $c_t : A \rightarrow [0, 1]$

## Definition

An *oblivious* adversary (or non-adaptive) is an adversary that depends only on  $A$  and  $t$ . Equivalently, an oblivious adversary has to choose the sequence of cost vectors at the beginning of time, knowing only the algorithm  $A$ .

In this lecture, we will only focus on the oblivious adversary.

# No regret Algorithm<sup>5</sup>

## Definition

If  $A$  is an online learning algorithm, then its expected regret at time  $t$  with respect to an arbitrary action  $a$  is

$$\mathbb{E} \left[ R_T^A(a) \right] = \frac{1}{T} \left( \sum_{i=1}^T \mathbb{E}_{a_t \sim \pi_t} [c_t(a_t)] - \sum_{t=1}^T c_t(a) \right),$$

where  $\{\pi_1, \pi_2, \dots, \pi_t\}$  is produced by  $A$

## Definition

An online decision-making algorithm  $A$  has no regret if for every  $\epsilon > 0$  there exists a sufficiently large time horizon  $T = T(\epsilon)$  such that, for every adversary for  $A$ , in expectation over the action realizations, the regret is at most  $\epsilon$ .

---

<sup>5</sup>We think of the number  $n$  of actions as fixed, and the time horizon  $T$  tending to infinity.

# Why do we care about online learning?

Consider a repeated normal-form game and suppose we are player  $i$  in the game. At time  $t \in \{1, 2, \dots, T\}$ :

- ▶ Each player  $j \neq i$  chooses some mixed strategy  $\pi_t^j$  over their pure strategies  $S^j$ .
- ▶ Let  $\sigma^{-i}$  be the product distribution over  $S^{-i}$  defined by these mixed strategies.
- ▶ Let  $c_t^i(a) = \mathbb{E}_{s^{-i} \sim \sigma^{-i}} [c^i(s^{-i}, a)]$

Then we managed to convert a repeated normal-form game to an online problem we just defined above.



# Why do we care about no regret algorithms?

Suppose you know nothing about your opponents in a game:

- ▶ By using a no-regret algorithm, player  $i$  guarantees that they do at least as well as the best action in hindsight.
- ▶ Since no-regret algorithms exist, why wouldn't players want to have this kind of guarantee?

What if you know something about your opponents in the game?

- ▶ You still have no-regret!

Therefore, any rational agent should have no-regret. It's a "minimum bar" for rationality!

# Connection to Equilibria

Suppose that we have a cost-minimisation game with  $k$  players and cost functions  $C^i : S \rightarrow \mathbb{R}$  for each player  $i \in [k]$ . Let's start with a few definitions:

- ▶ Let  $\pi_t^i$  be the mixed strategy used by player  $i$  at time  $t$ .
- ▶ Let  $\sigma^t = \prod_{i=1}^k \pi_t^i$  be the product distribution over  $S$  defined by the individual player distributions.
- ▶ Let  $\sigma = \frac{1}{T} \sum_{t=1}^T \sigma^t$  be the “average” distribution<sup>6</sup>.
- ▶ For each player  $i$  and time  $t$ , we define the cost vector  $c_t^i(a) = \mathbb{E}_{s \sim \sigma_t} [C^i(s^{-i}, a)]$ .

---

<sup>6</sup>Note that  $\sigma$  is not a product distribution since the players are correlated through  $t$ .

# Connection to Equilibria

Suppose that player  $i$  uses some algorithm  $A^i$  to generate its mixed strategy at each time.

## Theorem

*Suppose that  $E \left[ R_T^{A^i}(a) \right] \leq \epsilon$  for every  $i \in [k]$  and  $a \in S^i$ . Then  $\sigma$  is an  $\epsilon$ -approximate coarse correlated equilibrium:*

$$E_{s \sim \sigma} [C^i(s)] \leq E_{s \sim \sigma} [C^i(s^{-i}, s'')] + \epsilon,$$

*for all  $i \in [k]$  and  $s'' \in S$ .*



# No-Regret Algorithm: Multiplicative Weights

---

**Algorithm 1** Multiplicative Weights

---

- 1: Initialize  $w_1(a) = 1$  for all  $a \in \mathbb{A}$ .
  - 2: For  $t = 1, 2, \dots, T$ :
  - 3: **repeat**
  - 4:     Play an action according to the distribution:  $\pi_t = \frac{w_t}{\sum_{a \in A} w_t(a)}$
  - 5:     Given the cost vector  $c_t$ , update the weights by setting:  
       $w^{t+1}(a) = w^t(a) \cdot (1 - \eta)^{c^t(a)}$ , for all  $a \in A$
  - 6: **until** convergence
  - 7: output:  $\pi = \frac{1}{T} \sum_t \pi_t$
- 

- ▶ if  $\eta$  is small then we're *exploring* the space. If  $\eta$  is big, focus on actions which have been good in the past - we're *exploiting*.
- ▶ Also known as *Randomized Weighted Majority* and as *Hedge*.

# Regret Analysis

Is Multiplicative Weights a no-regret algorithm (expected external regret is at most  $2\sqrt{\frac{\ln |\mathcal{A}|}{T}}$ )?

**Proof.**

<sup>7</sup> Exercise (plenty of literature online).



---

<sup>7</sup>Typically we need to set  $\eta$  to  $\sqrt{\frac{\ln |\mathcal{A}|}{T}}$ .

# Correlated Equilibria: New Definition

## Definition

Let  $\sigma$  be a distribution over  $S = S^1 \times \dots \times S^k$ . Then  $\sigma$  is a correlated equilibrium if

$$\mathbb{E}_{s \sim \sigma} [c^i(s) \mid s^i] \leq \mathbb{E}_{s \sim \sigma} [c^i(s^{-i}, s^{i'}) \mid s^i],$$

for all  $i \in [k]$  and for all  $s^{i'} \in S^i$ .

This implies that player  $i$  does not want to switch from  $s^i$  to  $s^{i'}$  when told to play  $s^i$ . This way of thinking gives a different definition, in terms of “switching functions”.

## Definition

A switching function is a function  $\delta : S_i \rightarrow S_i$ , which we think of as telling to “switch” from some action  $s^i$  to  $\delta(s^i)$ .

# Correlated Equilibria: New Definition

## Definition

Let  $\sigma$  be a distribution over  $S = S^1 \times \dots \times S^k$ . Then  $\sigma$  is a correlated equilibrium if

$$E_{s \sim \sigma} [C^i(s)] \leq E_{s \sim \sigma} [C^i(s^{-i}, \delta(s^i))] ,$$

for all  $i \in [k]$  and for all  $\delta : S_i \rightarrow S_i$ .



# Correlated Equilibria and Swap Regret

This new definition of correlated equilibria will be useful because it will naturally relate back to a concept in online learning theory known as *swap regret*.

Recall we have defined our regret with respect to an action  $a$  as  $R_T(a) = \frac{1}{T} \left( \sum_{t=1}^T c_t(a_t) - \sum_{t=1}^T c_t(a) \right)$ . We call it external regret to distinguish from the swap regret defined below:

## Definition

The swap regret of a sequence of actions  $\{a_1, a_2, \dots, a_t\}$  with respect to a switching function  $\delta : A \rightarrow A$  is

$$S_T(\delta) = \frac{1}{T} \left( \sum_{t=1}^T c_t(a_t) - \sum_{t=1}^T c_t(\delta(a_t)) \right)$$

# Correlated Equilibria and Swap Regret

## Definition

The swap regret of a sequence of actions  $\{a_1, a_2, \dots, a_t\}$  with respect to a switching function  $\delta : A \rightarrow A$  is

$$S_T(\delta) = \frac{1}{T} \left( \sum_{t=1}^T c_t(a_t) - \sum_{t=1}^T c_t(\delta(a_t)) \right)$$

- ▶ Note that swap regret generalises our old notion of regret, since  $R_T(a) = S_T(\delta)$  where  $\delta(x) \equiv a$  for all  $x \in \mathbb{A}$ .
- ▶ So if we have low swap regret (for all  $\delta$ ) then we definitely have low regret (for all  $a$ ),
- ▶ but if we have low regret we might still have high swap regret.
- ▶ That is we're competing against a broader set of things in hindsight, so low swap regret is a stronger guarantee than low regret.

# More Definitions

## Definition

Let  $A$  be an online learning algorithm. Then its expected swap regret with respect to  $\delta : A \rightarrow A$  is

$$\mathbb{E} \left[ S_T^A(\delta) \right] = \frac{1}{T} \left( \sum_{t=1}^T \mathbb{E}_{a_t \sim \pi_t} [c_t(a_t)] - \sum_{t=1}^T \mathbb{E}_{a^t \sim \pi^t} [c_t(\delta(a_t))] \right)$$

## Definition

An online decision-making algorithm  $A$  has no swap regret if for every  $\epsilon > 0$  there exists a sufficiently large time horizon  $T = T(\epsilon)$  such that, for every adversary for  $A$ , the expected swap regret is at most  $\epsilon$ .

## Correlated Equilibria and Swap Regret (same as the previous setting)

Suppose that we have a cost-minimisation game with  $k$  players and cost functions  $C^i : S \rightarrow \mathbb{R}$  for each player  $i \in [k]$ . Let's start with a few definitions:

- ▶ Let  $\pi_t^i$  be the mixed strategy used by player  $i$  at time  $t$ .
- ▶ Let  $\sigma^t = \prod_{i=1}^k \pi_t^i$  be the product distribution over  $S$  defined by the individual player distributions.
- ▶ Let  $\sigma = \frac{1}{T} \sum_{t=1}^T \sigma^t$  be the “average” distribution<sup>8</sup>.
- ▶ For each player  $i$  and time  $t$ , we define the cost vector  $c_t^i(a) = \mathbb{E}_{s \sim \sigma_t} [C^i(s^{-i}, a)]$ .

---

<sup>8</sup>Note that  $\sigma$  is not a product distribution since the players are correlated through  $t$

# Correlated Equilibria and Swap Regret

## Theorem

Suppose that  $E \left[ S_T^{A^i}(\delta) \right] \leq \epsilon$  for all  $i \in [k]$  and for all  $\delta : S_i \rightarrow S_i$ .  
Then  $\sigma$  is an  $\epsilon$ -approximate correlated equilibrium, in the sense that

$$E_{s \sim \sigma} [C^i(s)] \leq E_{s \sim \sigma} [C^i(s^{-i}, \delta(s^i))] + \epsilon,$$

for all  $i \in [k]$  and for all  $\delta : S_i \rightarrow S_i$ .

# Correlated Equilibria and Swap Regret

Proof.

Let  $i \in [k]$  and  $\delta : S_i \rightarrow S_j$ . Recall that the cost vector that player  $i$  sees at time  $t$  is  $c_t(a) = \mathbb{E}_{s \sim \sigma_t} [C^i(s^{-i}, a)]$ , and hence  $\mathbb{E}_{a_t \sim \pi_t} [c_t(a)] = \mathbb{E}_{s \sim \sigma_t} [C^i(s)]$ . Then we have that:

$$\begin{aligned} & \mathbb{E}_{s \sim \sigma} [C^i(s)] - \mathbb{E}_{s \sim \sigma} [C^i(s^{-i}, \delta(s^i))] \\ &= \frac{1}{T} \sum_{i=1}^T \mathbb{E}_{s \sim \sigma^t} [C^i(s)] - \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{s \sim \sigma^t} [C^i(s^{-i}, \delta(s^i))] \\ &= \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{a_t \sim \pi_t} [c_t(a_t)] - \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{a_t \sim \pi_t} [c_t(\delta(a_t))] \\ &= \mathbb{E} [S_T^{A^i}(\delta)] \\ &\leq \epsilon \end{aligned}$$



