

Lecture 5

Learning in Games: Algorithms

Instructor: Gabriele Farina*

In Lecture 4 we have mentioned how no-external-regret dynamics recover several solution concepts of interest, including Nash equilibria in two-player zero-sum games, normal-form coarse-correlated equilibria in multiplayer general-sum games, and more generally convex-concave saddle point problems. Furthermore, we have seen that no-external-regret dynamics are a fundamental building block for constructing no- Φ -regret dynamics for general sets Φ of deviation functions.

In this lecture, we start exploring how no-external-regret dynamics for players in normal-form games can be constructed. As a reminder, a player with n action has a strategy space that corresponds to the probability simplex of all distributions over those actions,

$$\Delta^m = \{(x_1, \dots, x_m) \in \mathbb{R}_{\geq 0}^m : x_1 + \dots + x_m = 1\}.$$

Furthermore, remember that constructing a regret minimizer for Δ^m means that we need to build a mathematical object that supports two operations:

- `NextStrategy()` has the effect that the regret minimizer will output an element $\mathbf{x}^{(t)} \in \Delta^m$;
- `ObserveUtility($u^{(t)}$)` provides the environment's feedback to the regret minimizer, in the form of a utility function $u^{(t)}$. For today, we will focus on a case of particular relevance for normal-form games: the case of *linear* utilities (as such are the utilities in a normal-form game). In particular, to fix notation, we will let

$$u^{(t)} : \mathbf{x} \mapsto (\mathbf{g}^{(t)})^\top \mathbf{x},$$

where $\mathbf{g}^{(t)} \in \mathbb{R}^n$ is the vector that defines the linear function (called with the letter \mathbf{g} to suggest the idea of it being the “*gradient vector*”). Since the utility in the game cannot be unbounded, we assume that the $\mathbf{g}^{(t)}$ can be arbitrary but *bounded in norm*.

In this notation, the (external) *regret* is defined as the quantity

$$\text{Reg}^{(T)} := \max_{\hat{\mathbf{x}} \in \Delta^m} \left\{ \sum_{t=1}^T \left(\langle \mathbf{g}^{(t)}, \hat{\mathbf{x}} \rangle - \langle \mathbf{g}^{(t)}, \mathbf{x}^{(t)} \rangle \right) \right\}.$$

Our goal is to make sure that the regret grows sublinearly in T no matter the utility vectors $\mathbf{g}^{(t)}$ chosen by the environment.

Most algorithms known today for the task operate by *prioritizing actions based on how much regret they have incurred*. In particular, a key quantity to define modern algorithms is the vector of cumulated actions regrets

$$\mathbf{r}^{(t)} := \sum_{\tau=1}^t \mathbf{g}^{(\tau)} - \langle \mathbf{g}^{(\tau)}, \mathbf{x}^{(\tau)} \rangle \mathbf{1}.$$

(Note that $\text{Reg}^{(t)} = \max_{i \in \{1, \dots, m\}} \mathbf{r}^{(t)}[i]$.) The natural question is: how to prioritize?

*MIT EECS. ✉ gfarina@mit.edu.

1 An algorithm that does NOT work: Follow-the-Leader

Perhaps the most natural idea would be to always play the action with the highest cumulated regret (breaking ties, say, lexicographically). After all, this is the action we wish the most we had played in the past. Unfortunately, this idea is known not to work. To see that, consider the following sequence of gradient vectors:

$$\mathbf{g}^{(1)} = \begin{pmatrix} 0 \\ 1/2 \end{pmatrix}, \quad \mathbf{g}^{(2)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad \mathbf{g}^{(3)} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad \mathbf{g}^{(4)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad \mathbf{g}^{(5)} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad \text{etc.}$$

In this case, at time 1, we pick the first action and score a utility of 0. We then compute

$$\mathbf{r}^{(1)} = \begin{pmatrix} 0 \\ 1/2 \end{pmatrix}.$$

and at time $t = 2$ pick the second action since it has the highest regret. This results in a score of 0 and an updated regret vector

$$\mathbf{r}^{(2)} = \begin{pmatrix} 1 \\ 1/2 \end{pmatrix}$$

So, at time $t = 3$, we will pick the first action, resulting again in a score of 0 and an updated regret

$$\mathbf{r}^{(3)} = \begin{pmatrix} 1 \\ 3/2 \end{pmatrix},$$

and so on... Overall, it's easy to see that in all this jumping around, the regrets grow linearly.

2 The Regret Matching (RM) algorithm

The regret matching algorithm picks probabilities proportional to the “ReLU” of the cumulated regrets, that is,

$$\mathbf{x}^{(t+1)} \propto [\mathbf{r}^{(t)}]^+ \tag{1}$$

whenever $[\mathbf{r}^{(t)}] \neq \mathbf{0}$, and an arbitrary point otherwise.

The algorithm is presented in pseudocode in Algorithm 1.

Theorem 2.1. The regret matching algorithm (Algorithm 1) is an external regret minimizer, and satisfies the regret bound $R^T \leq \Omega\sqrt{T}$, where Ω is the maximum norm $\|\mathbf{g}^{(t)} - \langle \mathbf{g}^{(t)}, \mathbf{x}^{(t)} \rangle \mathbf{1}\|_2$ up to time T .

Proof. We start by observing that, at all times t ,

$$\left(\mathbf{g}^{(t+1)} - \langle \mathbf{g}^{(t+1)}, \mathbf{x}^{(t+1)} \rangle \mathbf{1} \right)^\top \mathbf{x}^{(t+1)} = 0$$

(this is always true, not just for regret matching). Plugging in the definition (1) of how $\mathbf{x}^{(t+1)}$ is constructed, we therefore conclude that

$$\left(\mathbf{g}^{(t+1)} - \langle \mathbf{g}^{(t+1)}, \mathbf{x}^{(t+1)} \rangle \mathbf{1} \right)^\top [\mathbf{r}^{(t)}]^+ = 0. \tag{2}$$

(Note that the above equation holds trivially when $[\mathbf{r}^{(t)}]^+$ and therefore $\mathbf{x}^{(t+1)}$ is picked arbitrarily.)

Now, we use the following inequality:

$$\|[\mathbf{a} + \mathbf{b}]^+\|_2^2 \leq \|[\mathbf{a}]^+ + \mathbf{b}\|_2^2,$$

applied to $\mathbf{a} = \mathbf{r}^{(t)}$ and $\mathbf{b} = \mathbf{g}^{(t+1)} - \langle \mathbf{g}^{(t+1)}, \mathbf{x}^{(t+1)} \rangle \mathbf{1}$. In particular, since by definition $\mathbf{a} + \mathbf{b} = \mathbf{r}^{(t+1)}$, we have

$$\begin{aligned} \|[\mathbf{r}^{(t+1)}]^+\|_2^2 &\leq \|[\mathbf{r}^{(t)}]^+ + (\mathbf{g}^{(t+1)} - \langle \mathbf{g}^{(t+1)}, \mathbf{x}^{(t+1)} \rangle \mathbf{1})\|_2^2 \\ &= \|[\mathbf{r}^{(t)}]^+\|_2^2 + \|\mathbf{g}^{(t+1)} - \langle \mathbf{g}^{(t+1)}, \mathbf{x}^{(t+1)} \rangle \mathbf{1}\|_2^2 + 2\left(\mathbf{g}^{(t+1)} - \langle \mathbf{g}^{(t+1)}, \mathbf{x}^{(t+1)} \rangle \mathbf{1}\right)^\top [\mathbf{r}^{(t)}]^+ \\ &= \|[\mathbf{r}^{(t)}]^+\|_2^2 + \|\mathbf{g}^{(t+1)} - \langle \mathbf{g}^{(t+1)}, \mathbf{x}^{(t+1)} \rangle \mathbf{1}\|_2^2 && \text{(from (2))} \\ &\leq \|[\mathbf{r}^{(t)}]^+\|_2^2 + \Omega^2. \end{aligned}$$

Hence, by induction we have

$$\|[\mathbf{r}^{(T)}]^+\|_2^2 \leq T\Omega^2 \implies \text{Reg}^{(T)} = \max_{i \in \{1, \dots, m\}} \mathbf{r}^{(T)}[i] \leq \max_{i \in \{1, \dots, m\}} [\mathbf{r}^{(T)}[i]]^+ \leq \|[\mathbf{r}^{(T)}]^+\|_2 \leq \Omega\sqrt{T}.$$

The proof is then complete. \square

Algorithm 1: Regret matching

```

1  $\mathbf{r}^0 \leftarrow \mathbf{0} \in \mathbb{R}^n, \mathbf{x}^0 \leftarrow \mathbf{1}/m \in \Delta^m$ 
2 function NextStrategy()
3   if  $[\mathbf{r}^{(t-1)}]^+ \neq \mathbf{0}$  return  $\mathbf{x}^{(t)} \leftarrow \frac{[\mathbf{r}^{(t-1)}]^+}{\|[\mathbf{r}^{(t-1)}]^+\|_1}$ 
4   else return  $\mathbf{x}^{(t)} \leftarrow$  any point in  $\Delta^m$ 
5 function ObserveUtility( $\mathbf{g}^t$ )
6    $\mathbf{r}^{(t)} \leftarrow \mathbf{r}^{(t-1)} + \mathbf{g}^{(t)} - \langle \mathbf{g}^{(t)}, \mathbf{x}^{(t)} \rangle \mathbf{1}$ 

```

Algorithm 2: Regret matching⁺

```

1  $\mathbf{z}^0 \leftarrow \mathbf{0} \in \mathbb{R}^n, \mathbf{x}^0 \leftarrow \mathbf{1}/m \in \Delta^m$ 
2 function NextStrategy()
3   if  $[\mathbf{r}^{(t-1)}]^+ \neq \mathbf{0}$  return  $\mathbf{x}^{(t)} \leftarrow \frac{[\mathbf{r}^{(t-1)}]^+}{\|[\mathbf{r}^{(t-1)}]^+\|_1}$ 
4   else return  $\mathbf{x}^{(t)} \leftarrow$  any point in  $\Delta^m$ 
5 function ObserveUtility( $\mathbf{g}^t$ )
6    $\mathbf{z}^{(t)} \leftarrow [\mathbf{z}^{(t-1)} + \mathbf{g}^{(t)} - \langle \mathbf{g}^{(t)}, \mathbf{x}^{(t)} \rangle \mathbf{1}]^+$ 

```

As of today, regret matching and its variants are still often some of the most practical algorithms for learning in games.

Remark 2.1. One very appealing property of the regret matching algorithm is its *lack of hyperparameters*. It just works “out of the box”.

2.1 The regret matching⁺ (RM⁺) algorithm

The regret matching⁺ algorithm was introduced by Tammelin [2014], Tammelin et al. [2015], and is given in Algorithm 2. It differs from RM only on the last line, where a further thresholding is added. That small change has the effect that actions with negative cumulated regret (that is, “bad” actions) are treated as actions with 0 regret. Hence, intuitively, if a bad action were to become good over time, it would take less time for RM⁺ to notice and act on that change. Because of that, regret matching⁺ has stronger practical performance and is often preferred over regret matching in the game solving literature.

With a simple modification to the analysis of RM, the same bound as RM can be proven.

Theorem 2.2. The RM^+ algorithm (Algorithm 2) is an external regret minimizer, and satisfies the regret bound $R^T \leq \Omega\sqrt{T}$, where Ω is the maximum norm $\|\mathbf{g}^{(t)} - \langle \mathbf{g}^{(t)}, \mathbf{x}^{(t)} \rangle \mathbf{1}\|_2$ up to time T .

3 Follow-the-Regularized-Leader (FTRL)

Another way of obtaining no-regret algorithms is by considering a regularized (i.e., smoothed) version of the follow-the-leader algorithm discussed above. The idea is that, instead of playing by always putting 100% of the probability mass on the action with highest cumulated regret, we look for the distribution that maximizes the expected cumulated regret, *minus* some regularization term that prevents us from putting all the mass on a single action.

More specifically, let $\varphi : \Delta^m \rightarrow \mathbb{R}$ be a 1-strongly convex¹ function (with respect to some norm $\|\cdot\|$, for example, the Euclidean norm $\|\cdot\|_2$) with bounded range, and consider the choice of strategy

$$\mathbf{x}^{(t)} := \arg \max_{\hat{\mathbf{x}} \in \Delta^m} \left\{ (\mathbf{r}^{(t-1)})^\top \hat{\mathbf{x}} - \frac{1}{\eta} \varphi(\hat{\mathbf{x}}) \right\}.$$

(so, in particular $\mathbf{x}^{(1)} := \arg \min_{\hat{\mathbf{x}} \in \Delta^m} \varphi(\hat{\mathbf{x}})$ since the regrets of all actions are 0 at the beginning). This algorithm is called the *follow-the-regularized-leader (FTRL)* algorithm. The following two properties are known about FTRL:

- Regularization *limits the amount of variation between consecutive strategies*. This makes intuitive sense: if $\eta \rightarrow 0$, for example, $\mathbf{x}^{(t)}$ is constant (and equal to $\arg \min_{\hat{\mathbf{x}} \in \Delta^m} \varphi(\hat{\mathbf{x}})$). When $\eta = \infty$, we recover the follow-the-leader algorithm, where strategies can jump arbitrarily. For intermediate η , as you might expect, the amount of variation between strategies at consecutive times is bounded above by a quantity proportional to η :

$$\|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\| \leq \eta \|\mathbf{g}^{(t)}\|_*, \quad (3)$$

where $\|\cdot\|_*$ is the *dual* norm of $\|\cdot\|$ (if $\|\cdot\|$ is the Euclidean norm, its dual is also the Euclidean norm; if $\|\cdot\|$ is the ℓ_1 norm, then its dual is the ℓ_∞ norm).

- A useful regret bound of FTRL is given by

$$\text{Reg}^{(T)} \leq \max_{\hat{\mathbf{x}} \in \Delta^m} \frac{\varphi(\hat{\mathbf{x}}) - \varphi(\mathbf{x}^{(1)})}{\eta} + \sum_{t=1}^T (\mathbf{g}^{(t)})^\top (\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}). \quad (4)$$

Using the Cauchy-Schwarz inequality in (4) and plugging in (3), we obtain the following

Theorem 3.1. The regret cumulated by the FTRL algorithm is upper bounded by

$$\text{Reg}^{(T)} \leq \max_{\hat{\mathbf{x}} \in \Delta^m} \frac{\varphi(\hat{\mathbf{x}}) - \varphi(\mathbf{x}^{(1)})}{\eta} + \eta \sum_{t=1}^T \|\mathbf{g}^{(t)}\|_*^2.$$

Since the norm of the gradient vectors is bounded, it follows immediately that by picking learning rate $\eta = 1/\sqrt{T}$, we obtain that $\text{Reg}^{(T)}$ is bounded as roughly \sqrt{T} (a sublinear function!) at all times T .

One of the appeals of FTRL is that it is a very general method, and in fact it generalized well beyond probability simplexes.

We will now see a particular instantiation of FTRL that is so important and well-studied that it deserves its own name: the *Multiplicative Weights Update* method (MWU).

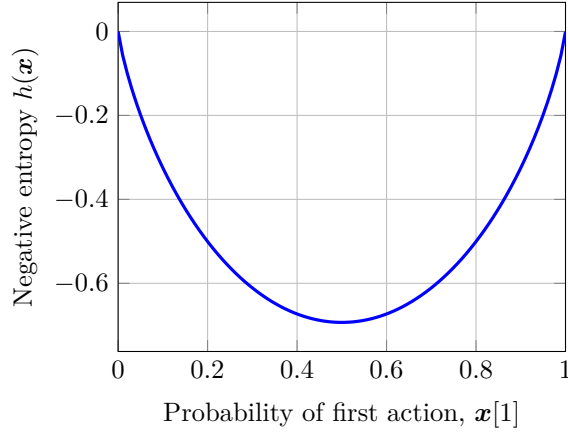
¹A differentiable function f is 1-strongly convex if for all \mathbf{x}, \mathbf{y} in its domain one has $(\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^\top (\mathbf{x} - \mathbf{y}) \geq \|\mathbf{x} - \mathbf{y}\|^2$.

3.1 The Multiplicative Weights Update Algorithm

A natural choice of convex regularizer for probability simplexes is the *negative entropy* function

$$h(\mathbf{x}) := \sum_{i \in \{1, \dots, m\}} \mathbf{x}[i] \log \mathbf{x}[i].$$

Example 3.1. For example, this is what negative entropy looks like in the case of $m = 2$ actions:



The negative entropy function has the following properties:

- It is 1-strongly convex with respect to the ℓ_1 norm $\|\cdot\|_1$.
- The maximum of the function is 0, which is attained at any deterministic strategy.
- Its minimum (that is, maximum entropy) is attained at the uniformly random strategy $\bar{\mathbf{x}} = (1/m, \dots, 1/m)$, at which the function has value $h(\bar{\mathbf{x}}) = m \cdot 1/m \log 1/m = -\log m$.

In this case, then, the MWU algorithm reduces to the following:

$$\begin{aligned} \mathbf{x}^{(1)} &:= (1/m, \dots, 1/m), \\ \mathbf{x}^{(t)} &:= \arg \max_{\hat{\mathbf{x}} \in \Delta^m} \left\{ (\mathbf{r}^{(t-1)})^\top \hat{\mathbf{x}} - \frac{1}{\eta} h(\hat{\mathbf{x}}) \right\} = \arg \max_{\hat{\mathbf{x}} \in \Delta^m} \left\{ (\mathbf{r}^{(t-1)})^\top \hat{\mathbf{x}} - \frac{1}{\eta} \sum_{i=1}^m \mathbf{x}[i] \log \mathbf{x}[i] \right\}. \end{aligned} \quad (5)$$

As it turns out via a simple application of the Lagrange multiplier theorem, the concave optimization problem in (5) has a closed-form solution:

Proposition 3.1. The concave optimization problem (5) has the closed-form solution

$$\mathbf{x}^{(t)} = \text{softmax}\{\eta \mathbf{r}^{(t)}\},$$

that is, for every action $i \in \{1, \dots, m\}$,

$$\mathbf{x}^{(t)}[i] = \frac{e^{\eta \mathbf{r}^{(t)}[i]}}{\sum_{j=1}^m e^{\eta \mathbf{r}^{(t)}[j]}}.$$

Algorithm 3: Multiplicative Weights Update

```
1  $\mathbf{r}^0 \leftarrow \mathbf{0} \in \mathbb{R}^n, \mathbf{x}^0 \leftarrow \mathbf{1}/m \in \Delta^m$ 

---

2 function NextStrategy()  
3   | return  $\mathbf{x}^{(t)} = \text{softmax}(\eta \mathbf{r}^{(t-1)})$ 

---

4 function ObserveUtility( $\mathbf{g}^t$ )  
5   |  $\mathbf{r}^{(t)} \leftarrow \mathbf{r}^{(t-1)} + \mathbf{g}^{(t)} - \langle \mathbf{g}^{(t)}, \mathbf{x}^{(t)} \rangle \mathbf{1}$ 

---


```

This gives rise to Algorithm 3, which has a similar flavor as RM, but uses a different prioritization of actions (in particular, it uses *softmax instead of ReLU*). From the general analysis of FTRL (Theorem 3.1), remembering that:

- the maximum of h is 0,
- the minimum (attained in $\mathbf{x}^{(1)}$) is $-\log m$,
- negative entropy is 1-strongly convex with respect to $\|\cdot\|_1$ and that the dual norm to $\|\cdot\|_1$ is $\|\cdot\|_\infty$,

we obtain the following.

Corollary 3.1. The regret cumulated by the MWU algorithm is upper bounded by

$$\text{Reg}^{(T)} \leq \frac{\log m}{\eta} + \eta \sum_{t=1}^T \|\mathbf{g}^{(t)}\|_\infty^2.$$

References

- Oskari Tammelin. Solving large imperfect information games using CFR+, 2014.
- Oskari Tammelin, Neil Burch, Michael Johanson, and Michael Bowling. Solving heads-up limit Texas hold'em. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, 2015.