# On the Optimality of Dilated Entropy and Lower Bounds for Online Learning in Extensive-Form Games

### Main question

Distance-generating functions provide generalized notions of distances Constrained optimization require a DGF to move appropriately in the feasible set

### What is the optimal DGF for extensive-form games?

### Key results

- Developed a new analysis for OMD based on the diameter/convexity ratio results of the DilEnt DGF.
- Established a matching regret lower bound, confirming the implied optimality of DilEnt.
- Achieved a new state-of-the-art convergence rate for Clairvoyant OMD to CCE.

Regularizer	Norm pair	$ \mathcal{D} /\mu$ ratio	Max gradient norm
Dilated Entropy [Kroer et al., 2020]	$\ell_1$ and $\ell_\infty$ norms	$\mathcal{O}(2^D \ \mathcal{Q}\ _1^2 \log  \mathcal{A} )$	$\leq 1$
Dilated Gl. Entropy [Farina et al., 2021]	$\ell_1$ and $\ell_\infty$ norms	$\mathcal{O}(\ \mathcal{Q}\ _1^2 \log  \mathcal{A} )$	$\leq 1$
DilEnt (this paper)	treeplex norms	$\ln  \mathcal{V} $	$\leq 1$

Table 1. Comparison of diameter/convexity ( $|\mathcal{D}|/\mu$ ) ratio with prior results. It holds that  $\ln |\mathcal{V}| \leq \mathcal{O}(\|\mathcal{Q}\|_1 \log |\mathcal{A}|)$ 

# **Tree-Form Sequential Decision Problem (TFSDP)**

**TFSDP** models the decision-making process of players in extensive-form games.

- A decision point  $j \in \mathcal{J}$  corresponds to an information set for the player.
- An observation point  $ja \in \Sigma$  represents the state that follows an action.



Figure 1. An example of an EFG and its corresponding TFSDP. Notably,  $p_{\rm B} = A1$ .

Sequence-form strategy space  $Q \subseteq \mathbb{R}^{\Sigma}$  structures strategies with linearity. Each entry x[ja] represents the probability of "reaching" the intermediate state ja, satisfying

$$\mathbf{x}[\varnothing] = 1, \qquad \sum_{a \in \mathcal{A}_j} \mathbf{x}[ja] = \mathbf{x}[p_j] \quad \forall j \in \mathcal{J}.$$

Here,  $p_j$  denote the preceding observation point of decision point j.

Zhiyuan Fan<sup>1</sup> Christian Kroer<sup>2</sup> Gabriele Farina<sup>1</sup>

<sup>2</sup>Columbia University <sup>1</sup>MIT

# **Problem Setting**

**Online learning** on sequence-form strategy space Q with full-information feedback:

- In each round, the agent picks  $\mathbf{x}_t \in \mathcal{Q}$ .
- They observe the adversarial reward vector  $\mathbf{w}_t$  and the reward  $\langle \mathbf{x}_t, \mathbf{w}_t \rangle \in [-1, 1]$ .

**Goal:** Minimize the cumulative regret:

 $\operatorname{Regret}(T) := \max_{\mathbf{x}_* \in \mathcal{Q}} \sum_{t=1}^{t} \langle \mathbf{x}_* - \mathbf{x}_t, \mathbf{w}_t \rangle.$ 

# **Proximal Method**

**Proximal step** generalizes the notion of gradient ascent to restricted space Q:

 $\mathbf{x}_{t+1} \leftarrow \Pi_{\varphi}(\eta \mathbf{g}, \mathbf{x}_t) := \operatorname*{argmax}_{\widehat{\mathbf{x}} \in \mathcal{O}} \{ \langle \eta \mathbf{g}, \widehat{\mathbf{x}} \rangle - \mathcal{D}_{\varphi}(\widehat{\mathbf{x}} \| \mathbf{x}_t) \}.$ 

Here,  $\varphi$  is a **distance-generating function (DGF)** and  $\mathcal{D}_{\varphi}$  is the Bregman divergence.

Given the primal-dual norms  $\|\cdot\|$  and  $\|\cdot\|_*$  with  $\|\mathbf{w}\|_* \leq 1$ , the performance of OMD is determined by:

- The strong convexity  $\mu$  of DGF over the primal norm  $\|\cdot$  $\triangleright$  A larger convexity  $\Rightarrow$  enables a larger step size  $\Rightarrow$  easier to learn
- The diameter  $|\mathcal{D}|$  of the strategy space measured by DGF ightarrow A smaller diameter  $\Rightarrow$  a smaller search space  $\Rightarrow$  easier to learn
- $\Rightarrow$  The diameter/convexity ratio  $|\mathcal{D}|/\mu$  is the key factor in the performance of mirror descent.

We inspect the weight-one dilated entropy (DilEnt) in TFSDP:

$$\varphi_{1}: \mathcal{Q} \ni \mathbf{x} \mapsto \sum_{j \in \mathcal{J}} \sum_{a \in \mathcal{A}} \mathbf{x}[ja] \ln \left( \frac{\mathbf{x}[ja]}{\mathbf{x}[p_{j}]} \right)$$

 $\triangleright$  When TFSDP only has one decision point,  $\varphi_1$  reduces to entropy.

# **Primal-Dual Treeplex Norms**

We introduce a pair of primal-dual norms: the treeplex  $\ell_1$ -norm  $\|\cdot\|_{\mathcal{H},1}$  and the treeplex  $\ell_{\infty}$ -norm  $\|\cdot\|_{\mathcal{H},\infty}$ :

$$\|\mathbf{u}\|_{\mathcal{H},1} := \sup_{\mathbf{y}\in\mathcal{Q}^*} \langle |\mathbf{u}|, \mathbf{y} \rangle,$$

Here,  $Q^*$  is the dual polytope of Q. Both treeplex norms can be computed via recursion:

$$\|\mathbf{u}\|_{\mathcal{H}_p,1} := \max_{j:p_j=p} \sum_{a \in \mathcal{A}_j} \|\mathbf{u}\|_{\mathcal{H}_{ja},1},$$

	•
	•



$$\|\mathbf{u}\|_{\mathcal{H},\infty} := \sup_{\mathbf{x}\in\mathcal{Q}} \langle |\mathbf{u}|, \mathbf{x} \rangle.$$

$$\|\mathbf{u}\|_{\mathcal{H}_p,\infty} := \sum_{j:p_j=p} \max_{a \in \mathcal{A}_j} \|\mathbf{u}\|_{\mathcal{H}_{ja},\infty}$$

[Theorem 5.6] If every player runs Clairvoyant OMD with DilEnt in an *n*-player game, the average joint policy converges to a *Coarse Correlated Equilibrium* at the rate

[Theorem 6.1 + Theorem 6.2] Any algorithm incurs an expected regret of at least

Under certain structural assumptions, the bound is free from hidden logarithmic factors.

DilEnt achieves an nearly-optimal diameter/convexity ratio:

functions. Mathematical Programming, pages 1–33, 2020.

# **Key Structural Results**

**[Lemma 4.4]** The dual norm of any feasible reward vector w satisfies  $\|w\|_{\mathcal{H},\infty} \leq 1$ .

[Lemma 5.1] DilEnt is 1-strongly convex with respect to the treeplex  $\ell_1$ -norm:  $\|\cdot\|^2_{\nabla^2 \omega_1(\mathbf{x})} \ge \|\cdot\|^2_{\mathcal{H},1}$ .

[Lemma 5.2] The diameter of the strategy space is upper bounded as  $\max_{x_*} \mathcal{D}_{\varphi_1}(x_*, x_1) \leq \ln |\mathcal{V}|$ .

### Main Results on Regret Upper Bounds

[Theorem 5.4] Running OMD with DilEnt during the proximal steps achieves a regret bound of

 $\operatorname{Regret}(T) \leq \sqrt{2\ln|\mathcal{V}|}\sqrt{T},$ 

where  $|\mathcal{V}| := \mathcal{Q} \cap \{0,1\}^{\Sigma}$  is the number of pure strategies.

 $\varepsilon \leq \mathcal{O}(n \log |\mathcal{V}| \log T/T).$ 

# Lower Bound Implies Optimality of DilEnt

 $\operatorname{Regret}(T) \ge \widetilde{\Omega}(\sqrt{\ln |\mathcal{V}|}\sqrt{T})$ 

• Better DGF  $\Rightarrow$  better regret upper bound in Theorem 5.4

• A contradiction arises since superior regret upper bounds violate the lower bound above.

Checkout our paper on openreview!



Gabriele Farina, Christian Kroer, and Tuomas Sandholm. Better regularization for sequential decision spaces: Fast convergence rates for nash, correlated, and team equilibria. In Proceedings of the 2021 ACM Conference on Economics and Computation, 2021.

Christian Kroer, Kevin Waugh, Fatma Kılınç-Karzan, and Tuomas Sandholm. Faster algorithms for extensive-form game solving via improved smoothing