

# Online Learning in Complex Environments: The Need for a Data-dependent Theory

**Tim van Erven**



UNIVERSITY  
OF AMSTERDAM

Game AI Workshop @ Maastricht

March 30, 2023

# Online Convex Optimization

Parameters  $\theta$  take values in a convex domain  $\Theta \subset \mathbb{R}^d$

- 1: **for**  $t = 1, 2, \dots, T$  **do**
- 2:   Learner estimates  $\theta_t \in \Theta$
- 3:   Nature reveals convex loss function  $f_t : \Theta \rightarrow \mathbb{R}$
- 4: **end for**

**Goal:** Predict almost as well as the best possible parameters  $\theta^*$ :

$$\text{Regret}_T(\theta^*) = \sum_{t=1}^T f_t(\theta_t) - \sum_{t=1}^T f_t(\theta^*)$$

# Online Convex Optimization

Parameters  $\theta$  take values in a convex domain  $\Theta \subset \mathbb{R}^d$

- 1: **for**  $t = 1, 2, \dots, T$  **do**
- 2:   Learner estimates  $\theta_t \in \Theta$
- 3:   Nature reveals convex loss function  $f_t : \Theta \rightarrow \mathbb{R}$
- 4: **end for**

**Goal:** Predict almost as well as the best possible parameters  $\theta^*$ :

$$\text{Regret}_T(\theta^*) = \sum_{t=1}^T f_t(\theta_t) - \sum_{t=1}^T f_t(\theta^*)$$

Viewed as a **zero-sum game** against Nature:

$$V = \min_{\theta_1} \max_{f_1} \min_{\theta_2} \max_{f_2} \cdots \min_{\theta_T} \max_{f_T} \max_{\theta^* \in \Theta} \text{Regret}_T(\theta^*)$$

# Example: Electricity Forecasting



Typically, functions  $f_t$  determined by data:

- ▶ Every day  $t$  an electricity company needs to predict how much electricity  $Y_t$  is needed the next day
- ▶ Given feature vector  $\mathbf{X}_t \in \mathbb{R}^d$ , predict  $\hat{Y}_t = \mathbf{X}_t^\top \boldsymbol{\theta}_t$  with a linear model
- ▶ Next day: observe  $Y_t$
- ▶ Measure loss by  $f_t(\boldsymbol{\theta}_t) = (Y_t - \hat{Y}_t)^2$  and improve parameter estimates:  $\boldsymbol{\theta}_t \rightarrow \boldsymbol{\theta}_{t+1}$

# Online Gradient Descent

$$\begin{aligned}\tilde{\theta}_{t+1} &= \theta_t - \eta_t \nabla f_t(\theta_t) \\ \theta_{t+1} &= \min_{\theta \in \Theta} \|\tilde{\theta}_{t+1} - \theta\|\end{aligned}$$

## Theorem (Zinkevich, 2003)

*Suppose  $\Theta$  compact with diameter at most  $D$ , and  $\|\nabla f_t(\theta_t)\| \leq G$ . Then online gradient descent with  $\eta_t = \frac{D}{G\sqrt{t}}$  guarantees*

$$\text{Regret}_T(\theta^*) \leq \frac{3}{2} GD\sqrt{T}$$

*for **any** choices of Nature.*

Without further assumptions, this is optimal (up to a constant factor).

# OGD is Optimal, but is it Good?

## Theorem (Lower Bound)

*For any learning algorithm, there exists an OCO task with  $\text{diam}(\Theta) \leq D$  and  $\|\nabla f_t(\theta_t)\| \leq G$  such that*

$$\text{Regret}_T(\theta^*) \geq cGD\sqrt{T}$$

*for some absolute constant  $c > 0$ .*

### Proof:

- ▶  $\Theta = [-\frac{D}{2}, \frac{D}{2}]$
- ▶  $f_t(\theta) = \theta g_t$  with  $\Pr(g_t = -G) = \Pr(g_t = +G) = 1/2$

Then for any algorithm

$$\mathbb{E} \left[ \sum_{t=1}^T f_t(\theta_t) \right] = 0,$$

but

$$\mathbb{E} \left[ \min_{\theta^* \in \Theta} \sum_{t=1}^T f_t(\theta^*) \right] = \frac{D}{2} \mathbb{E} \left[ \min \left\{ \sum_{t=1}^T g_t, -\sum_{t=1}^T g_t \right\} \right] \leq -cDG\sqrt{T}.$$

# OGD is Optimal, but is it Good?

## Theorem (Lower Bound)

*For any learning algorithm, there exists an OCO task with  $\text{diam}(\Theta) \leq D$  and  $\|\nabla f_t(\theta_t)\| \leq G$  such that*

$$\text{Regret}_T(\theta^*) \geq cGD\sqrt{T}$$

*for some absolute constant  $c > 0$ .*

### Proof:

- ▶  $\Theta = [-\frac{D}{2}, \frac{D}{2}]$
- ▶  $f_t(\theta) = \theta g_t$  with  $\Pr(g_t = -G) = \Pr(g_t = +G) = 1/2$

### Hardest case:

- ▶ **Linear functions**  $f_t$
- ▶ **Gradients**  $g_t$  are **pure noise**, with maximal variance  
→ nothing interesting to learn, so irrelevant for applications

# What if there is Less Noise?

## Theorem (Sachs, Hadiji, Van Erven, Guzmán, 2023)

*There exists an algorithm (optimistic follow-the-regularized-leader) that guarantees the worst-case bound*

$$\text{Regret}_T(\boldsymbol{\theta}^*) = O(GD\sqrt{T})$$

*and, if the  $f_t$  are i.i.d. and  $F_t(\boldsymbol{\theta}) = \mathbb{E}[f_t(\boldsymbol{\theta})]$  is  $L$ -smooth, then*

$$\mathbb{E}[\text{Regret}_T(\boldsymbol{\theta}^*)] = O(\sigma D\sqrt{T} + LD^2)$$

- ▶  $\sigma^2 = \max_{\boldsymbol{\theta} \in \Theta} \text{Var}(\nabla f_t(\boldsymbol{\theta}))$
- ▶ Exploits stochasticity and smoothness if available, but does **not assume** them
- ▶ Previously known for linear losses, i.e.  $L = 0$  [Rakhlin, Sridharan, 2013]
- ▶ Recovers optimal rates in stochastic acceleration (via online-to-batch conversion on scaled losses)



# What if there is Less Noise? (Refined Version)

## Theorem (Sachs, Hadiji, Van Erven, Guzmán, 2023)

*There exists an algorithm (optimistic follow-the-regularized-leader) that guarantees the worst-case bound*

$$\text{Regret}_T(\theta^*) = O(GD\sqrt{T})$$

*and, if the  $f_t$  are stochastic and each  $F_t(\theta) = \mathbb{E}[f_t(\theta)|\mathcal{F}_{t-1}]$  is  $L$ -smooth, then*

$$\mathbb{E}[\text{Regret}_T(\theta^*)] = O((\bar{\sigma}_T + \bar{\Sigma}_T)D\sqrt{T} + LD^2)$$

- ▶  $\bar{\sigma}_T^2 = \frac{1}{T} \sum_{t=1}^T \max_{\theta \in \Theta} \text{Var}(\nabla f_t(\theta))$ : average **variance** of the gradients
- ▶  $\bar{\Sigma}_T^2 = \frac{1}{T} \sum_{t=1}^T \max_{\theta \in \Theta} \|\nabla F_t(\theta) - \nabla F_{t-1}(\theta)\|^2$ : average **drift** in the expected gradients
- ▶ **Interpolates** between i.i.d. and adversarial settings

# Towards a Data-dependent Theory of Online Learning

## Applications are not zero-sum games:

1. Worst-case regret witnessed on **fully random data**
  - ▶ Not relevant for practice!
2. **Nature is not trying to win**: e.g.
  - ▶ Consumers do not adjust electricity consumption to make statistical analysis hard

Can often adapt to some data- or distribution-dependent measure of **easiness of the data** to get much better performance!

# Standard Textbook View of General OCO [Hazan, 2016]

Convex $f_t$	$\sqrt{T}$	Online Gradient Descent with $\eta_t \propto \frac{1}{\sqrt{t}}$
Strongly convex $f_t$	$\ln T$	Online Gradient Descent with $\eta_t \propto \frac{1}{t}$
Exp-concave $f_t$	$d \ln T$	Online Newton Step with $\eta \propto 1$

## Minimax rates based on curvature (bounded domain and gradients)

- **Strongly convex:** second derivative at least  $\alpha > 0$ , implies exp-concave
- **Exp-concave:**  $e^{-\alpha \ell_t}$  concave  
Satisfied by log loss, logistic loss, squared loss, but not hinge loss

# Standard Textbook View of General OCO [Hazan, 2016]

Convex $f_t$	$\sqrt{T}$	Online Gradient Descent with $\eta_t \propto \frac{1}{\sqrt{t}}$
Strongly convex $f_t$	$\ln T$	Online Gradient Descent with $\eta_t \propto \frac{1}{t}$
Exp-concave $f_t$	$d \ln T$	Online Newton Step with $\eta \propto 1$

**Minimax rates based on curvature** (bounded domain and gradients)

## Limitations:

- ▶ Different method in each case. (Requires sophisticated users.)
- ▶ Theoretical tuning of  $\eta_t$  **very conservative**
- ▶ What if curvature varies between rounds?
- ▶ In many applications data are **stochastic** (i.i.d.) Should be easier than worst case. . .

# Standard Textbook View of General OCO [Hazan, 2016]

Convex $f_t$	$\sqrt{T}$	Online Gradient Descent with $\eta_t \propto \frac{1}{\sqrt{t}}$
Strongly convex $f_t$	$\ln T$	Online Gradient Descent with $\eta_t \propto \frac{1}{t}$
Exp-concave $f_t$	$d \ln T$	Online Newton Step with $\eta \propto 1$

**Minimax rates based on curvature** (bounded domain and gradients)

## Limitations:

- ▶ Different method in each case. (Requires sophisticated users.)
- ▶ Theoretical tuning of  $\eta_t$  **very conservative**
- ▶ What if curvature varies between rounds?
- ▶ In many applications data are **stochastic** (i.i.d.) Should be easier than worst case. . .

## Need Adaptive Methods!

- ▶ Difficulty: All existing methods learn  $\eta$  at too slow rate [HP2005] so **overhead of learning best  $\eta$  ruins potential benefits**

## Theorem (Van Erven, Koolen, 2016, Van Erven, Koolen, Van der Hoeven, 2021)

The MetaGrad algorithm guarantees the following **data-dependent** bound:

$$\text{Regret}_T(\theta^*) \leq \sum_{t=1}^T (\theta_t - \theta^*)^\top \nabla f_t(\theta_t) \preceq \begin{cases} \sqrt{T \ln \ln T} \\ \sqrt{V_T(\theta^*) d \ln T} + d \ln T \end{cases}$$

where

$$V_T(\theta^*) = \sum_{t=1}^T ((\theta^* - \theta_t)^\top \nabla f_t(\theta_t))^2.$$

### Key Feature:

- Pay only  $\ln \ln T$  for learning  $\eta$

# Consequences

## 1. Non-stochastic adaptation:

Convex $f_t$	$\sqrt{T \ln \ln T}$
Exp-concave $f_t$	$d \ln T$
Fixed convex $f_t = f$	$d \ln T$

Extension by [Wang, Lu, Zhang, 2020] also achieves  $O(\ln T)$  for strongly convex losses

# Consequences

## 1. Non-stochastic adaptation:

Convex $f_t$	$\sqrt{T \ln \ln T}$
Exp-concave $f_t$	$d \ln T$
Fixed convex $f_t = f$	$d \ln T$

Extension by [Wang, Lu, Zhang, 2020] also achieves  $O(\ln T)$  for strongly convex losses

## 2. Stochastic without curvature [Koolen, Grünwald, Van Erven, 2016]:

Suppose  $f_t$  i.i.d. with stochastic optimum  $\theta^* = \arg \min_{\theta \in \Theta} \mathbb{E}_f[f(\theta)]$ .

Then expected regret  $\mathbb{E}[\text{Regret}_T(\theta^*)]$ :

Absolute loss* $f_t(\theta) =  \theta - X_t $ for $d=1$	$\ln T$
Hinge loss* $\max\{0, 1 - Y_t \langle \theta, X_t \rangle\}$	$d \ln T$
<b><math>(B, \beta)</math>-Bernstein</b>	$(Bd \ln T)^{1/(2-\beta)} T^{(1-\beta)/(2-\beta)}$

\*Conditions apply



# MetaGrad Experiments [Van Erven, Koolen, Van der Hoeven, 2021]

- ▶ 17 benchmark **UCI data sets**: 11 classification, 6 regression
- ▶ Tune algorithms according to theory  
(**No secret hyperparameter optimization!**)
- ▶ Measure regret ratio between algorithm and Online Gradient Descent

Algorithm	Median Regret Ratio	Computation Time
AdaGrad	3.54	$O(dT)$
Online Gradient Descent	1.00	$O(dT)$
MetaGrad	<b>0.25</b>	$O(d^2 T)$

# MetaGrad Experiments [Van Erven, Koolen, Van der Hoeven, 2021]

- ▶ 17 benchmark **UCI data sets**: 11 classification, 6 regression
- ▶ Tune algorithms according to theory  
(**No secret hyperparameter optimization!**)
- ▶ Measure regret ratio between algorithm and Online Gradient Descent

Algorithm	Median Regret Ratio	Computation Time
AdaGrad	3.54	$O(dT)$
Online Gradient Descent	1.00	$O(dT)$
MetaGrad	<b>0.25</b>	$O(d^2 T)$

But... MetaGrad **slow in high dimensions**. Fast approximations:

Coordinatewise	0.32	$O(dT)$
Sketching( $m = 1$ )	0.31	$O(mdT)$
Sketching( $m = 10$ )	0.27	$O(mdT)$
Sketching( $m = 50$ )	0.25	$O(mdT)$

# Many Possible Sources of Easiness in Data

- ▶ Statistical: nature is (approximately) stationary
- ▶ Curvature: loss function is benign
- ▶ Game-theoretic: other players update slowly + smoothness
- ▶ Model selection: maybe a simple comparator  $\theta^*$  is optimal
- ▶ Structured comparator classes
- ▶ Smoothed analysis: data have smooth distribution
- ▶ Bandits: less exploration needed
- ▶ ...?

# Food for Thought

## Can organize by:

- ▶ Application domain: games, bandits, full information, market making, optimization, bilateral trade, ...
- ▶ Types of adaptivity: types of loss functions, easiness caused by statistical or deterministic regularity, adapting to hyperparameters like  $G$  or  $\|\theta^*\|$ , ...
- ▶ Techniques: adaptive learning rates, optimistic gradient estimates, adaptive exploration for bandits, ...

## Working group goals:

- ▶ List desirable types of adaptivity for various settings
- ▶ Organize them
- ▶ Prioritize
- ▶ Identify expertise, common interests and collaborations

# Food for Thought

## Can organize by:

- ▶ Application domain: games, bandits, full information, market making, optimization, bilateral trade, ...
- ▶ Types of adaptivity: types of loss functions, easiness caused by statistical or deterministic regularity, adapting to hyperparameters like  $G$  or  $\|\theta^*\|$ , ...
- ▶ Techniques: adaptive learning rates, optimistic gradient estimates, adaptive exploration for bandits, ...

## Working group goals:

- ▶ List desirable types of adaptivity for various settings
- ▶ Organize them
- ▶ Prioritize
- ▶ Identify expertise, common interests and collaborations

## No silver bullet:

- ▶ The price of adaptivity: if overhead (computational/regret) too large, may not be worth it
- ▶ Some types of adaptivity may be mutually exclusive, e.g.  $G$  vs  $\|\theta^*\|$ , or may be impossible in some settings.

# References

- E. Hazan. Introduction to online convex optimization. *Foundations and Trends in Optimization*, 2(3–4):157–325, 2016.
- W. M. Koolen, P. Grünwald, and T. van Erven. Combining adversarial guarantees and stochastic fast rates in online learning. In *Advances in Neural Information Processing Systems 29 (NeurIPS)*, pages 4457–4465, 2016.
- A. Rakhlin and K. Sridharan. Online learning with predictable sequences. In *Proceedings of the 26th Annual Conference on Learning Theory (COLT)*, pages 993–1019, 2013.
- T. van Erven and W. M. Koolen. Metagrad: Multiple learning rates in online learning. In *Advances in Neural Information Processing Systems 29 (NeurIPS)*, pages 3666–3674, 2016.
- T. van Erven, W. M. Koolen, and D. van der Hoeven. Metagrad: Adaptation using multiple learning rates in online learning. *Journal of Machine Learning Research*, 22(161):1–61, 2021. URL <http://jmlr.org/papers/v22/20-1444.html>.