

Online Convex Optimization and Its Surprising Applications

Francesco Orabona

KAUST

Learning Theory Summer School, Copenhagen, Denmark, 2026



جامعة الملك عبد الله
للعلوم والتقنية
King Abdullah University of
Science and Technology

Parameter-free Online Algorithms

What is a Parameter-free Algorithm?

Definition

We call parameter-free any online algorithm that satisfies the optimal regret bound with respect to T and to all the comparators in the feasible set \mathcal{V} , up to poly-logarithmic multiplicative factors.

- Exponentiated Gradient: $Regret_T(\mathbf{u}) \leq \frac{KL(\mathbf{u}; \pi)}{\eta} + \frac{T\eta}{2} \Rightarrow$
NormalHedge: $Regret_T(\mathbf{u}) = \mathcal{O}(\sqrt{T(KL(\mathbf{u}; \pi) + 1)})$ [Chaudhuri et al., NeurIPS'09][Chernov&Vovk, UAI'10][Orabona&Pál, NeurIPS'16]
- OSD: $Regret_T(\mathbf{u}) \leq \frac{\|\mathbf{x}_1 - \mathbf{u}\|_2^2}{2\eta} + \frac{\eta T}{2} \Rightarrow$
KT (next slides): $Regret_T(\mathbf{u}) = \mathcal{O}(\|\mathbf{x}_1 - \mathbf{u}\|_2 \sqrt{T \ln(1 + T\|\mathbf{x}_1 - \mathbf{u}\|_2/\epsilon)} + \epsilon)$
- How do I tune my learning rate/regularizer without knowing \mathbf{u} ?

Parameter-Free through Duality on Guarantee

- We saw the online-to-batch conversion (deterministic case for simplicity):

$$F(\bar{\mathbf{x}}_T) - F(\mathbf{u}) \leq \frac{1}{T} \sum_{t=1}^T (F(\mathbf{x}_t) - F(\mathbf{u})) \leq \frac{1}{T} \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{u} \rangle \leq \frac{\psi_T(\mathbf{u})}{T}$$

Theorem (McMahan&Orabona, COLT'14)

Consider two sequences $\mathbf{x}_1, \dots, \mathbf{x}_T$ and $\mathbf{g}_1, \dots, \mathbf{g}_T$. Then

$$\begin{aligned} \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{u} \rangle &\leq \psi_T(\mathbf{u}), \quad \forall \mathbf{u} \in \mathbb{R}^d \\ &\iff \\ -\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t \rangle &\geq \psi_T^* \left(-\sum_{t=1}^T \mathbf{g}_t \right), \quad \forall \mathbf{g}_1, \dots, \mathbf{g}_T \end{aligned}$$

where ψ_T^* is the Fenchel conjugate of ψ_T defined as $\psi_T^*(\boldsymbol{\theta}) = \sup_{\mathbf{x}} \langle \boldsymbol{\theta}, \mathbf{x} \rangle - \psi_T(\mathbf{x})$

- Assume $\|\mathbf{g}_t\|_2 \leq 1$
- Set $\mathbf{x}_t = \frac{-\sum_{i=1}^{t-1} \mathbf{g}_i}{t} \left(1 - \sum_{i=1}^{t-1} \langle \mathbf{g}_i, \mathbf{x}_i \rangle\right)$
- Claim: \mathbf{x}_t guarantees

$$-\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t \rangle \geq \psi_T^* \left(-\sum_{t=1}^T \mathbf{g}_t \right)$$

where $\psi_T^*(\boldsymbol{\theta}) \approx \frac{1}{\sqrt{T}} \exp\left(\frac{\|\boldsymbol{\theta}\|_2^2}{2T}\right) - 1$

- This implies $\sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{u} \rangle \leq \|\mathbf{x}^*\|_2 \sqrt{T \ln(\|\mathbf{u}\|_2 T + 1)} + 1$
- Where do the update rule and the inequality in orange come from?

Optimization through Optimal Gambling

Krichevsky&Trofimov (KT) betting strategy:

- Observe sequence of coins outcomes $c_t \in [-1, 1]$, start with \$1, bet x_t money, win/lose $x_t c_t$
- On round t bet a signed fraction of your money equal to $\frac{\sum_{i=1}^{t-1} c_i}{t}$
- Exponential amount of money

$$\text{Wealth of KT after } T \text{ rounds} = 1 + \sum_{t=1}^T x_t c_t \geq \frac{\exp\left(\frac{(\sum_{t=1}^T c_t)^2}{2T}\right)}{2\sqrt{T}}$$

- No assumptions on the coin!
- We need to prove that $-\sum_{t=1}^T g_t x_t \geq \psi_T^* \left(-\sum_{t=1}^T g_t\right)$
- In 1d, set $c_t = -g_t$ and assume $|g_t| \leq 1$ then we have it!
- It works in the vector case too

- It works in any number of dimensions, even Hilbert spaces

$$\mathbf{x}_t = \mathbf{x}_0 + \frac{-\sum_{i=1}^{t-1} \mathbf{g}_i}{t} \left(1 - \sum_{i=1}^{t-1} \langle \mathbf{g}_i, \mathbf{x}_i \rangle \right)$$

- It works with stochastic subgradients
- Variant that does not need to know the Lipschitz constant [Cutkosky, COLT'19]
- Variant works with constrained sets [Cutkosky&Orabona, COLT'18]
- Variant adapts to the strong convexity in the stochastic setting (bounded stochastic subgradients and domain) [Cutkosky&Orabona, COLT'18]

Surprising Applications of Online Learning

Online Learning is Much More than Online Learning

- Online Convex Optimization might seem only concerned with losses®ret
- In reality, it is about proving inequalities on arbitrary sequences of data
- Sometimes, the inequalities are more important than the algorithms

- Here, I'll try to convince you of this view

From Online Convex Optimization to Non-Convex Non-Smooth Optimization

a.k.a.

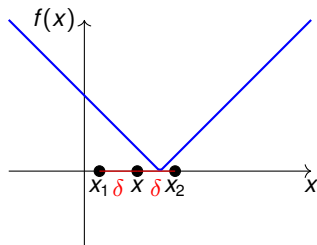
Online losses can be anything

Non-convex Optimization

- For convex optimization, we study $F(\mathbf{x}_T) - F(\mathbf{u})$
- For non-convex smooth optimization, we study $\mathbb{E}_i[\|\nabla F(\mathbf{x}_i)\|_2^2]$
- What can we do for non-smooth non-convex? Example: ConvNets with ReLUs

Definition (Zhang et al. ICML'20)

A point \mathbf{x} is an (δ, ϵ) -stationary point of an almost-everywhere differentiable function F if there is a finite subset S of the ball of radius δ centered at \mathbf{x} such that for \mathbf{y} selected uniformly at random from S , $\mathbb{E}[\mathbf{y}] = \mathbf{x}$ and $\|\mathbb{E}[\nabla F(\mathbf{y})]\| \leq \epsilon$



If δ is small enough, it codifies our intuition on points close to a minimum

We will assume that the functions are well-behaved in the sense that

$$F(\mathbf{y}) - F(\mathbf{x}) = \int_0^1 \langle \nabla F(\mathbf{x} + t(\mathbf{y} - \mathbf{x})), \mathbf{y} - \mathbf{x} \rangle dt$$

Up to perturbing the function with some noise, this holds for locally Lipschitz functions

Using OCO for Non-convex Optimization

Require: An OCO algorithm, duration of cycle K , initial point \mathbf{x}_0

- 1: $j = 0$
- 2: **for** $t = 1$ **to** T **do**
- 3: **if** $\text{mod}(t, K) == 1$ **then**
- 4: Reset OCO algorithm
- 5: $j = j + 1$
- 6: $\bar{\mathbf{x}}_j = \mathbf{0}$
- 7: **end if**
- 8: Receive \mathbf{m}_t from OCO algorithm
- 9: $\mathbf{x}_t = \mathbf{x}_{t-1} - \mathbf{m}_t$
- 10: Sample s_t uniformly in $[0, 1]$
- 11: $\mathbf{x}'_t = \mathbf{x}_{t-1} - s_t \mathbf{m}_t$
- 12: Pass $\ell_t(\mathbf{x}) = -\langle \nabla F(\mathbf{x}'_t), \mathbf{x} \rangle$ to OCO algorithm
- 13: $\bar{\mathbf{x}}_j = \bar{\mathbf{x}}_j + \mathbf{x}'_t / K$
- 14: **end for**
- 15: **return** $\bar{\mathbf{x}}_J$ where J is drawn uniformly at random between 1 and T/K

■ **Important:** The OCO algorithm decides the updates not the iterates

Theorem

Let the OCO algorithm be OGD over the L_2 ball of radius D . Then, we have

$$\mathbb{E} \left[\frac{1}{T/K} \sum_{i=1}^{T/K} \left\| \frac{1}{K} \sum_{t=1}^K \nabla F(\mathbf{x}'_{(i-1)K+t}) \right\|_2 \right] \leq \frac{F(\mathbf{x}_0) - \inf_{\mathbf{x}} F(\mathbf{x})}{DT} + \frac{1}{\sqrt{K}}$$

Moreover, set $D = \delta/K$, $K = \left(\frac{T\delta}{F(\mathbf{x}_0) - \inf_{\mathbf{x}} F(\mathbf{x})} \right)^{\frac{2}{3}}$, and return $\bar{\mathbf{x}}_J$ where J is uniformly at random, then in expectation $\bar{\mathbf{x}}_J$ is $(\delta, \mathcal{O}((T\delta)^{-\frac{1}{3}}))$ -stationary point.

- The choice of D : $\bar{\mathbf{x}}_j$ is the average of K points contained in ball of radius δ
- With the chosen D , we have

$$\frac{F(\mathbf{x}_0) - \inf_{\mathbf{x}} F(\mathbf{x})}{DT} + \frac{1}{\sqrt{K}} = \frac{K(F(\mathbf{x}_0) - \inf_{\mathbf{x}} F(\mathbf{x}))}{T\delta} + \frac{1}{\sqrt{K}}$$

- $\mathbb{E} \left[\frac{1}{T/K} \sum_{i=1}^{T/K} \left\| \frac{1}{K} \sum_{t=1}^K \nabla F(\mathbf{x}'_{(i-1)K+t}) \right\|_2 \right] =$
 $\mathbb{E} \left[\left\| \frac{1}{K} \sum_{t=1}^K \nabla F(\mathbf{x}'_{(J-1)K+t}) \right\|_2 \right] = \mathcal{O} \left((T\delta)^{-\frac{1}{3}} \right)$

[Cutkosky et al., ICML'23]

In all optimization analyses we need to link function values to gradients:

- Convex functions: $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$
- Non-convex M -smooth: $f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{M}{2} \|\mathbf{x} - \mathbf{y}\|^2$
- What can we use for non-convex non-smooth?

Key Observation

- We evaluate the gradient in $\mathbf{x}'_t = \mathbf{x}_{t-1} - s_t \mathbf{m}_t = \mathbf{x}_{t-1} + s_t(\mathbf{x}_t - \mathbf{x}_{t-1})$
- Hence, we have

$$\mathbb{E}_{s_t} \nabla F(\mathbf{x}'_t) = \int_0^1 \nabla F(\mathbf{x}_{t-1} + t(\mathbf{x}_t - \mathbf{x}_{t-1})) dt$$

- This allows us to say that

$$\begin{aligned} F(\mathbf{x}_t) - F(\mathbf{x}_{t-1}) &= \int_0^1 \langle \nabla F(\mathbf{x}_{t-1} + t(\mathbf{x}_t - \mathbf{x}_{t-1})), \mathbf{x}_t - \mathbf{x}_{t-1} \rangle dt \\ &= \langle \mathbb{E}_{s_t} [\nabla F(\mathbf{x}'_t)], \mathbf{x}_t - \mathbf{x}_{t-1} \rangle \end{aligned}$$

- This holds without assuming convexity nor smoothness!

Proof.

Using the key observation, for the first cycle we have

$$\begin{aligned} F(\mathbf{x}_t) - F(\mathbf{x}_{t-1}) &= \langle \mathbb{E}_{S_t}[\nabla F(\mathbf{x}'_t)], \mathbf{x}_t - \mathbf{x}_{t-1} \rangle = -\langle \mathbb{E}_{S_t}[\nabla F(\mathbf{x}'_t)], \mathbf{m}_t \rangle \\ &= \langle -\mathbb{E}_{S_t}[\nabla F(\mathbf{x}'_t)], \mathbf{m}_t - \mathbf{u} \rangle - \langle \mathbb{E}_{S_t}[\nabla F(\mathbf{x}'_t)], \mathbf{u} \rangle \end{aligned}$$

Taking full expectation, summing over $t = 1, \dots, K$, for any $\|\mathbf{u}\|_2 \leq D$, we have

$$\begin{aligned} \mathbb{E}[F(\mathbf{x}_K)] - F(\mathbf{x}_0) &= \mathbb{E} \left[\underbrace{\sum_{t=1}^K \langle -\nabla F(\mathbf{x}'_t), \mathbf{m}_t - \mathbf{u} \rangle}_{\text{Regret}_K(\mathbf{u})} \right] - \mathbb{E} \left[\sum_{t=1}^K \langle \nabla F(\mathbf{x}'_t), \mathbf{u} \rangle \right] \\ &\leq D\sqrt{K} - \mathbb{E} \left[\sum_{t=1}^K \langle \nabla F(\mathbf{x}'_t), \mathbf{u} \rangle \right] \end{aligned}$$

Choose $\mathbf{u} = D \frac{\sum_{t=1}^K \nabla F(\mathbf{x}'_t)}{\left\| \sum_{t=1}^K \nabla F(\mathbf{x}'_t) \right\|_2}$ to have $\sum_{t=1}^K \langle \nabla F(\mathbf{x}'_t), \mathbf{u} \rangle = -D \left\| \sum_{t=1}^K \nabla F(\mathbf{x}'_t) \right\|_2$.

Summing over the cycles and dividing by DT ends the proof. \square

Only a Hack or Something Fundamental?

One might wonder if the above reduction is only a hack or it discovers something more fundamental.

One way to convince you is to take a look at the resulting procedure

$$\begin{aligned}\mathbf{x}_t &= \mathbf{x}_{t-1} - \mathbf{m}_t \\ \mathbf{g}_t &= \nabla F(\mathbf{x}_t + (s_t - 1)\mathbf{m}_t) \\ \mathbf{m}_{t+1} &= \text{Clip}_D(\mathbf{m}_t + \eta\mathbf{g}_t)\end{aligned}$$

We recovered a version of SGD with momentum and clipping! The only really different part is that we perturb the iterate a bit before calculating the gradient

From Online Betting to PAC-Bayes Bounds

a.k.a.

The existence of online algorithms implies
generalization bounds

Definitions:

$$\text{Risk}(\theta) = \mathbb{E}_{(x,y) \sim P}[\ell(y, f_{\theta}(x))] \quad (\text{True risk of } \theta)$$

$$\text{Risk}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f_{\theta}(X_i)) \quad (\text{Empirical risk of } \theta)$$

Assumption:

$$0 \leq \ell \leq 1$$

Theorem (McAllester, COLT'98)

Fix a prior distribution $\pi \in \mathcal{M}(\Theta)$. With probability at least $1 - \delta$ on the data $S = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$, for any probability distribution ρ learnt on the data,

$$\mathbb{E}_{\theta \sim \rho}[\text{Risk}(\theta)] \leq \mathbb{E}_{\theta \sim \rho}[\text{Risk}_n(\theta)] + \sqrt{\frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{n}}{\delta}}{2n}}$$

Theorem

Define optimal 'log-wealth' function:

$$\psi_n^*(\theta) := \max_{\lambda \in [-\frac{1}{1 - \text{Risk}(\theta)}, \frac{1}{\text{Risk}(\theta)}]} \sum_{i=1}^n \ln[1 + \lambda(\ell(Y_i, f_\theta(X_i)) - \text{Risk}(\theta))] .$$

Fix $\pi \in \mathcal{M}(\Theta)$, then with probability at least $1 - \delta$, **simultaneously** for all n and ρ ,

$$\mathbb{E}_{\theta \sim \rho}[\psi_n^*(\theta)] \leq \text{KL}(\rho \parallel \pi) + \ln \frac{\sqrt{n}}{\delta} .$$

I. By relaxing the 'log-wealth' term this inequality implies:

- McAllester's inequality [McAllester, COLT'98]
- Empirical Bernstein's PAC-Bayes inequality [Tolstikhin&Seldin, NeurIPS'13]
- Maurer's inequality of Bernoulli r.v.'s [Maurer, arXiv'04]
- Unexpected Bernstein's inequality [Mhammedi et al., NeurIPS'19]

II. With no relaxations, we can compute confidence sequences on μ_θ **efficiently**

Our inequality:

$$\begin{aligned} \mathbb{E}_{\theta \sim \rho}[\psi_n^*(\theta)] &:= \mathbb{E}_{\theta \sim \rho} \left[\max_{\lambda \in \left[-\frac{1}{1 - \text{Risk}(\theta)}, \frac{1}{\text{Risk}(\theta)}\right]} \sum_{i=1}^n \ln(1 + \lambda(\ell(Y_i, f_\theta(X_i)) - \text{Risk}(\theta))) \right] \\ &\leq \text{KL}(\rho \parallel \pi) + \ln \frac{\sqrt{n}}{\delta} \end{aligned}$$

- $\ln(1 + x) \geq x - x^2$ for $x \geq -0.68$ gives

$$|\mathbb{E}_{\theta \sim \rho}[\text{Risk}(\theta)] - \mathbb{E}_{\theta \sim \rho}[\text{Risk}_n(\theta)]| \leq 2\sqrt{\frac{\text{KL}(\rho \parallel \pi) + \ln \frac{\sqrt{n}}{\delta}}{n}} \Rightarrow \text{McAllester's bound!}$$

- By convexity, $\max_{\lambda} \sum_{i=1}^n \ln(1 + \lambda(X_i - \mu)) \geq n \text{kl}(\hat{\mu}, \mu)$, that gives

$$\text{kl}\left(\mathbb{E}_{\theta \sim \rho}[\text{Risk}_n(\theta)], \mathbb{E}_{\theta \sim \rho}[\text{Risk}(\theta)]\right) \leq \frac{\text{KL}(\rho \parallel \pi) + \ln \frac{\sqrt{n}}{\delta}}{n} \Rightarrow \text{Maurer's bound!}$$

- Similarly, you can get the other bounds too

Let's Start from the Classic Proof

For any $\rho \ll \pi$ and measurable F :

$$\mathbb{E}_{\theta \sim \rho}[F(\theta)] \leq \text{KL}(\rho \parallel \pi) + \ln \mathbb{E}_{\theta \sim \pi}[e^{F(\theta)}] \quad (\text{Change-of-measure})$$

For some fixed $\lambda > 0$ choose $F(\theta) = \lambda(\text{Risk}(\theta) - \text{Risk}_n(\theta))$. Then,

$$\begin{aligned} \lambda \mathbb{E}_{\theta \sim \rho}[\text{Risk}(\theta) - \text{Risk}_n(\theta)] &\leq \text{KL}(\rho \parallel \pi) + \ln \mathbb{E}_{\theta \sim \pi}[e^{\lambda(\text{Risk}(\theta) - \text{Risk}_n(\theta))}] \\ &\leq \text{KL}(\rho \parallel \pi) + \ln \frac{1}{\delta} + \ln \mathbb{E}_{\theta \sim \pi}[\mathbb{E} e^{\lambda(\text{Risk}(\theta) - \text{Risk}_n(\theta))}] \end{aligned} \quad (\text{Markov})$$

Concentration!

E.g., Hoeffding's lemma + tuning over λ gives McAllester's inequality.

Standard approach: λ is fixed. **Idea:** tune λ based on data...

... using an online betting game:

- A fictitious betting algorithm starts with wealth 1
- At round $i = 1, \dots, n$ it bets a signed fraction of its wealth $B_i(\theta)$
- Observes outcome $\Delta_i(\theta) := \ell(Y_i, f_\theta(X_i)) - \text{Risk}(\theta)$
- Then its log wealth is $\psi_n(\theta) := \sum_{i=1}^n \ln(1 + B_i(\theta)\Delta_i(\theta))$
- The regret of the algorithm is controlled, as before:

$$\psi_n^*(\theta) - \psi_n(\theta) \leq \ln \sqrt{n}, \quad \forall \theta$$

where $\psi_n^*(\theta)$ is the log-wealth of the optimal constant betting fraction:

$$\psi_n^*(\theta) = \max_{\lambda \in \left[-\frac{1}{1 - \text{Risk}(\theta)}, \frac{1}{\text{Risk}(\theta)}\right]} \sum_{i=1}^n \ln[1 + \lambda(\ell(Y_i, f_\theta(X_i)) - \text{Risk}(\theta))]$$

For any $\rho \ll \pi$ and measurable F :

$$\mathbb{E}_{\theta \sim \rho}[F(\theta)] \leq \text{KL}(\rho \parallel \pi) + \ln \mathbb{E}_{\theta \sim \pi}[e^{F(\theta)}] \quad (\text{Change-of-measure})$$

Choose $F(\theta) = \psi_n^*(\theta)$ (optimal log-wealth). Then,

$$\mathbb{E}_{\theta \sim \rho}[\psi_n^*(\theta)] \leq \text{KL}(\rho \parallel \pi) + \ln \mathbb{E}_{\theta \sim \pi}[e^{\psi_n^*(\theta)}]$$

$$e^{\psi_n^*(\theta)} = \text{OptimalWealth} \leq \text{WealthAnyOnlineAlgorithm}_n \cdot \exp(\text{Regret}_n(A))$$

Concentration: $\text{WealthAnyOnlineAlgorithm}_n$ is a non-negative martingale

$$\Pr \left\{ \sup_{n \geq 0} \text{WealthAnyOnlineAlgorithm}_n \geq \frac{1}{\delta} \right\} \leq \delta \quad (\text{Ville's inequality})$$

Putting all together, with probability at least $1 - \delta$, we have

$$\mathbb{E}_{\theta \sim \rho}[\psi_n^*(\theta)] \leq \text{KL}(\rho \parallel \pi) + \ln \left(\frac{1}{\delta} \exp(\ln \sqrt{n}) \right) = \text{KL}(\rho \parallel \pi) + \ln \frac{1}{\delta} + \ln \sqrt{n}$$

Normalized Gradients for All

a.k.a.

You don't lose anything in using the adversarial
setting

Two Assumptions, Two Learning Rates

- Consider M -smooth functions, that is, $\|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})\|_* \leq M\|\mathbf{x} - \mathbf{y}\|$
- It is well-known that gradient descent converges as $\mathcal{O}(1/T)$ on smooth functions, with learning rate $\eta = \frac{1}{M}$
- On the other hand, we only have $\mathcal{O}(1/\sqrt{T})$ on Lipschitz functions with learning rate $\eta \propto \frac{1}{\sqrt{T}}$
- Two different assumptions need two different learning rates

- Can I have a single algorithm that gives both rates, without anything to tune?
- Is there anything in between these two rates?

- Let's generalize smoothness: $\|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})\|_* \leq M_\nu \|\mathbf{x} - \mathbf{y}\|^\nu$, for $\nu \in [0, 1]$
- $\nu = 0$ is for Lipschitz functions
- $\nu = 1$ is the usual smooth assumption

- For smooth function, we have $\|\nabla F(\mathbf{x})\|_*^2 \leq 2M(F(\mathbf{x}) - F(\mathbf{x}^*))$
- For Hölder-smooth functions, we have $\|\nabla F(\mathbf{x})\|_*^{1+\frac{1}{\nu}} \leq (1 + \frac{1}{\nu}) M_\nu^{\frac{1}{\nu}} (F(\mathbf{x}) - F(\mathbf{x}^*))$

- Very simple idea: take your favourite first-order optimization algorithm
- Use normalized gradients instead of gradients

- If your algorithm satisfies a regret guarantee on Lipschitz losses, then it automatically adapts to all the Hölder-smooth functions!

Theorem

Suppose that we have an algorithm that guarantees

$$\sum_{t=1}^T \langle \mathbf{q}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq R_T(\mathbf{x}^*)$$

for any sequence of $\mathbf{q}_1, \dots, \mathbf{q}_T$ such that $\|\mathbf{q}_t\|_* \leq 1$. Also, assume F is M_ν Hölder-smooth.

Then, the same algorithm fed with $\mathbf{q}_t = \frac{\nabla F(\mathbf{x}_t)}{\|\nabla F(\mathbf{x}_t)\|_*}$ guarantees

$$F(\bar{\mathbf{x}}_T) - F(\mathbf{x}^*) \leq 2M_\nu \left(\frac{R_T(\mathbf{x}^*)}{T} \right)^{1+\nu}$$

where $\bar{\mathbf{x}}_T = \frac{1}{\sum_{t=1}^T \frac{1}{\|\nabla F(\mathbf{x}_t)\|_*}} \sum_{t=1}^T \frac{\mathbf{x}_t}{\|\nabla F(\mathbf{x}_t)\|_*}$

- Consider Gradient Descent with $\eta = \frac{1}{\sqrt{T}}$, it satisfies

$$\sum_{t=1}^T \langle \mathbf{q}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \frac{\sqrt{T}}{2} \|\mathbf{x}_1 - \mathbf{x}^*\|_2^2 + \frac{\sqrt{T}}{2}$$

- How do we know it? Because we proved it yesterday!
- So, Gradient Descent with normalized gradients satisfies

$$F(\bar{\mathbf{x}}_T) - F(\mathbf{x}^*) \leq 2M_\nu \left(\frac{\|\mathbf{x}_1 - \mathbf{x}^*\|_2^2 + 1}{2\sqrt{T}} \right)^{1+\nu}$$

- Known since Levy [NeurIPS'17] for $\nu \in \{0, 1\}$
- Strictly more general than Mishchenko and Malitsky [ICML'20] that works only for $\nu = 1$
- Rediscovered in DoGW [Khaled et al. NeurIPS'23] for $\nu \in \{0, 1\}$

- Consider the KT algorithm I mentioned before, it satisfies that

$$\sum_{t=1}^T \langle \mathbf{q}_t, \mathbf{x}_t - \mathbf{x}^* \rangle \leq \mathcal{O}(\|\mathbf{x}_1 - \mathbf{x}^*\|_2 \sqrt{T \ln(T\|\mathbf{x}_1 - \mathbf{x}^*\|_2 + 1)})$$

- So, KT with normalized gradients satisfies

$$F(\bar{\mathbf{x}}_T) - F(\mathbf{x}^*) \leq \mathcal{O} \left(M_\nu \left(\frac{\|\mathbf{x}_1 - \mathbf{x}^*\|_2 \sqrt{T \ln(T\|\mathbf{x}_1 - \mathbf{x}^*\|_2 + 1)}}{\sqrt{T}} \right)^{1+\nu} \right)$$

- Known since Orabona&Pal [arXiv'21]
- Rediscovered in DoGW [Khaled et al. NeurIPS'23] for $\nu = 0$ and $\nu = 1$

$$\sum_{t=1}^T \frac{F(\mathbf{x}_t) - F(\mathbf{x}^*)}{\|\nabla F(\mathbf{x}_t)\|_*} \stackrel{\text{Convexity}}{\leq} \sum_{t=1}^T \left\langle \frac{\nabla F(\mathbf{x}_t)}{\|\nabla F(\mathbf{x}_t)\|_*}, \mathbf{x}_t - \mathbf{x}^* \right\rangle \stackrel{\text{Assumption}}{\leq} R_T(\mathbf{x}^*) \quad (*)$$

Using $\bar{\mathbf{x}}_T = \frac{1}{\sum_{t=1}^T \frac{1}{\|\nabla F(\mathbf{x}_t)\|_*}} \sum_{t=1}^T \frac{\mathbf{x}_t}{\|\nabla F(\mathbf{x}_t)\|_*}$, we have

$$\begin{aligned} F(\bar{\mathbf{x}}_T) - F(\mathbf{x}^*) &\stackrel{(*) + \text{Jensen}}{\leq} \frac{1}{\sum_{t=1}^T \frac{1}{\|\nabla F(\mathbf{x}_t)\|_*}} R_T(\mathbf{x}^*) \stackrel{\text{HM-GM}}{\leq} \frac{R_T(\mathbf{x}^*)}{T} \left(\prod_{t=1}^T \|\nabla F(\mathbf{x}_t)\|_* \right)^{\frac{1}{T}} \\ &= \frac{R_T(\mathbf{x}^*)}{T} \left(\prod_{t=1}^T \frac{\|\nabla F(\mathbf{x}_t)\|_*^{1+\frac{1}{\nu}}}{\|\nabla F(\mathbf{x}_t)\|_*} \right)^{\frac{\nu}{\nu+1}} \\ &\stackrel{\text{smoothness}}{\leq} \frac{R_T(\mathbf{x}^*)}{T} \left(\prod_{t=1}^T \frac{(1+\frac{1}{\nu}) M_\nu^{\frac{1}{\nu}} (F(\mathbf{x}_t) - F(\mathbf{x}^*))}{\|\nabla F(\mathbf{x}_t)\|_*} \right)^{\frac{\nu}{\nu+1}} \\ &\stackrel{\text{GM-AM}}{\leq} \frac{(1+\frac{1}{\nu})^\nu M_\nu R_T(\mathbf{x}^*)}{T} \left(\frac{1}{T} \sum_{t=1}^T (F(\mathbf{x}_t) - F(\mathbf{x}^*)) \right)^{\nu} \stackrel{(*)}{\leq} 2M_\nu \left(\frac{R_T(\mathbf{x}^*)}{T} \right)^{\nu+1} \end{aligned}$$

Even More Surprising Applications

- Rademacher complexity bounds from Online Learning [Kakade et al., NeurIPS'08]
- Better-than-KL PAC-Bayes bounds [Kuzborskij et al., COLT'24]
- Parameter-free sampling [Sharrock&Nemeth, ICML'23][Sharrock et al. NeurIPS'23]

- Confidence intervals through betting (my talk in the workshop)

- Basic concepts and definitions of Online Learning
- OMD&FTRL
- Parameter-free algorithms
- Connection between regret and betting, PAC-Bayes bounds, and adaptation

Thank you!

Website: <https://francesco.orabona.com>

Blog: <https://parameterfree.com>

X/Twitter: @bremen79