

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

Aims of the Lecture

- Provide an introduction to Online Convex Optimization
- *Almost* rigorous: details are missing, but theorems are correct
- (1-slide) Proofs! Because it is the only way to design online learning algorithms
- Ideally, when in 1 week all this material will disappear from your memory, you can still use the slides as a “cheat sheet”
- Most of the material is based on my upcoming book on online learning (<https://arxiv.org/abs/1912.13213>), my blog posts (<https://parameterfree.com>), and some recent papers

- 1 Online Convex Optimization and Regret
- 2 Online Mirror Descent
- 3 Follow-the-Regularized-Leader
- 4 Parameter-free Online Algorithms
- 5 From Online Learning to Non-smooth Non-convex Optimization
- 6 From Online Betting to PAC-Bayes
- 7 From Online Learning to Adaptation to Smoothness

Online Learning

- 1 In each round, output $\mathbf{x}_t \in \mathcal{V}$
- 2 Pay $\ell_t(\mathbf{x}_t)$
- 3 Update \mathbf{x}_{t+1} based on received information on ℓ_t

Choose \mathbf{x}_t before observing ℓ_t

No assumptions on how ℓ_t is generated!

Regret minimization

$$\min_{\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathcal{V}} \sum_{t=1}^T \ell_t(\mathbf{x}_t) \quad \text{equivalently} \quad \min_{\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathcal{V}} \underbrace{\sum_{t=1}^T \ell_t(\mathbf{x}_t) - \sum_{t=1}^T \ell_t(\mathbf{u})}_{\text{Regret}_T(\mathbf{u})}$$

- The algorithm is *no-regret* if $\lim_{T \rightarrow \infty} \frac{1}{T} \text{Regret}_T(\mathbf{u}) \leq 0$ for all $\mathbf{u} \in \mathcal{V}$ and any sequence of losses in a certain family

Why Online Convex Optimization?

- It is a strict generalization of the learning with expert setting
- It generalizes the setting of batch and stochastic convex optimization, in 99% of the cases without losing anything
- It provides a different mindset for designing optimization algorithms
- It is connected to a number of topics: Adaptation, Generalization, PAC-Bayes, Compression, Betting, etc.

Some Famous Online Learning Algorithms

- Online Gradient Descent [Zinkevich, ICML'03]
- AdaGrad [Duchi et al., COLT'10, JMLR'11; McMahan&Streeter, COLT'10]
- AMSGrad [Reddi et al., ICLR'18]

These algorithms are designed to work in the adversarial setting and have a $\mathcal{O}(\sqrt{T})$ regret bound

We will see that they can also be used as stochastic optimization algorithms with a $\mathcal{O}(\frac{1}{\sqrt{T}})$ convergence rate

- Losses: $\ell_t : \mathbb{R}^d \rightarrow \mathbb{R}$, convex, 1-Lipschitz
- Feasible set: $\mathcal{V} \subseteq \mathbb{R}^d$, closed, convex, non-empty
- Iterates: All technical conditions for iterates \mathbf{x}_t to exist hold

Mainly Two Main Meta-Algorithms

- Online Mirror Descent (OMD)
- Follow-the-Regularized-Leader (FTRL)

- These two meta-algorithms cover 90% of the (online) optimization algorithms
- Examples
 - Online Gradient Descent = special case of OMD
 - Dual Averaging = Special case of FTRL with linearized losses
 - Regularized Dual Averaging = Special case of FTRL with linearized losses
 - “Lazy version” of online gradient descent = FTRL
 - Newton algorithm = OMD with distance induced by the Hessian
 - Accelerated algorithm = two OCO algorithms playing against each other
 - Frank-Wolfe algorithm = two OCO algorithms playing against each other
 - etc.

Online Subgradient Descent

Require: Feasible set $\mathcal{V} \subseteq \mathbb{R}^d$, $\mathbf{x}_1 \in \mathcal{V}$, $\eta_1, \dots, \eta_T > 0$

1: **for** $t = 1$ **to** T **do**

2: Output $\mathbf{x}_t \in \mathcal{V}$

3: Pay $\ell_t(\mathbf{x}_t)$

4: Set $\mathbf{g}_t = \nabla \ell_t(\mathbf{x}_t)$

5: $\mathbf{x}_{t+1} = \Pi_{\mathcal{V}}(\mathbf{x}_t - \eta_t \mathbf{g}_t) = \operatorname{argmin}_{\mathbf{y} \in \mathcal{V}} \|\mathbf{x}_t - \eta_t \mathbf{g}_t - \mathbf{y}\|_2$

6: **end for**

Lemma

Let $\ell_t : \mathcal{V} \rightarrow \mathbb{R}$ differentiable in an open set that contains \mathcal{V} . Then, $\forall \mathbf{u} \in \mathcal{V}$, OGD satisfies

$$\eta_t(\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) \leq \eta_t \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{u} \rangle \leq \frac{1}{2} \|\mathbf{x}_t - \mathbf{u}\|_2^2 - \frac{1}{2} \|\mathbf{x}_{t+1} - \mathbf{u}\|_2^2 + \frac{\eta_t^2}{2} \|\mathbf{g}_t\|_2^2 .$$

Proof.

$$\begin{aligned} \|\mathbf{x}_{t+1} - \mathbf{u}\|_2^2 - \|\mathbf{x}_t - \mathbf{u}\|_2^2 &\stackrel{\text{\(\(\Pi\)\ is non expansive}}{\leq} \|\mathbf{x}_t - \eta_t \mathbf{g}_t - \mathbf{u}\|_2^2 - \|\mathbf{x}_t - \mathbf{u}\|_2^2 \\ &= -2\eta_t \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{u} \rangle + \eta_t^2 \|\mathbf{g}_t\|_2^2 \\ &\stackrel{\text{Convexity}}{\leq} -2\eta_t(\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) + \eta_t^2 \|\mathbf{g}_t\|_2^2 . \end{aligned}$$

□

Theorem

Let ℓ_1, \dots, ℓ_T differentiable in open sets containing \mathcal{V} . Pick any $\mathbf{x}_1 \in \mathcal{V}$ and assume $\eta_t = \eta$, $t = 1, \dots, T$. Then, $\forall \mathbf{u} \in \mathcal{V}$, OGD satisfies

$$\sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) \leq \frac{\|\mathbf{u} - \mathbf{x}_1\|_2^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{g}_t\|_2^2 - \frac{1}{2\eta} \|\mathbf{x}_{T+1} - \mathbf{u}\|_2^2.$$

Proof.

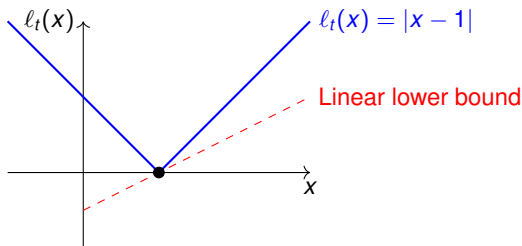
Dividing the inequality in the previous Lemma by η and summing over $t = 1, \dots, T$, we have

$$\begin{aligned} \sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) &\leq \sum_{t=1}^T \left(\frac{1}{2\eta} \|\mathbf{x}_t - \mathbf{u}\|_2^2 - \frac{1}{2\eta} \|\mathbf{x}_{t+1} - \mathbf{u}\|_2^2 \right) + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{g}_t\|_2^2 \\ &= \frac{1}{2\eta} \|\mathbf{x}_1 - \mathbf{u}\|_2^2 - \frac{1}{2\eta} \|\mathbf{x}_{T+1} - \mathbf{u}\|_2^2 + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{g}_t\|_2^2. \end{aligned}$$

□

Non-Differentiable Convex Functions

- If the losses are convex, but not differentiable, we cannot calculate the gradients
- We only need gradients because they satisfy
$$\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u}) \leq \langle \nabla \ell_t(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle$$
- Solution: use any vector \mathbf{g}_t that satisfies $\ell_t(\mathbf{x}_t) - \ell(\mathbf{u}) \leq \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{u} \rangle$ for all $\mathbf{u} \in \mathcal{V}$
- \mathbf{g}_t is called a **subgradient** of ℓ_t in \mathbf{x}_t
- The set of all subgradients ℓ in \mathbf{x} is called **subdifferential** and it is denoted by $\partial \ell_t(\mathbf{x}_t)$



Projected Online Subgradient Descent

Require: Feasible set $\mathcal{V} \subseteq \mathbb{R}^d$, $\mathbf{x}_1 \in \mathcal{V}$, $\eta_1, \dots, \eta_T > 0$

- 1: **for** $t = 1$ **to** T **do**
- 2: Output $\mathbf{x}_t \in \mathcal{V}$
- 3: Pay $\ell_t(\mathbf{x}_t)$
- 4: Set $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$
- 5: $\mathbf{x}_{t+1} = \Pi_{\mathcal{V}}(\mathbf{x}_t - \eta_t \mathbf{g}_t) = \operatorname{argmin}_{\mathbf{y} \in \mathcal{V}} \|\mathbf{x}_t - \eta_t \mathbf{g}_t - \mathbf{y}\|_2$
- 6: **end for**

Same guarantee of OGD:

$$\sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) \leq \frac{\|\mathbf{u} - \mathbf{x}_1\|_2^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{g}_t\|_2^2 - \frac{1}{2\eta} \|\mathbf{x}_{T+1} - \mathbf{u}\|_2^2.$$

- The regret is $\sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) \leq \frac{\|\mathbf{u} - \mathbf{x}_1\|_2^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{g}_t\|_2^2$
- Assume the function 1-Lipschitz w.r.t. the L_2 norm
($|\ell_t(\mathbf{x}) - \ell_t(\mathbf{u})| \leq \|\mathbf{x} - \mathbf{u}\|_2$)
- Then, $\sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) \leq \frac{\|\mathbf{u} - \mathbf{x}_1\|_2^2}{2\eta} + \frac{T\eta}{2}$
- Optimal learning rate: $\eta = \frac{\|\mathbf{u} - \mathbf{x}_1\|_2}{\sqrt{T}}$
- Any problem with this choice?

- Practical choice $\eta = \frac{\alpha}{\sqrt{T}}$ that gives $\text{Regret}_T(\mathbf{u}) \leq \frac{1}{2} \left(\frac{\|\mathbf{x}_1 - \mathbf{u}\|_2^2}{\alpha} + \alpha \right) \sqrt{T}$

- Easy case: \mathcal{V} has bounded diameter D , then $\eta = \frac{D}{\sqrt{T}}$ gives regret $D\sqrt{T}$

Applications: From Online to Stochastic (or Batch) Optimization (1)

- 1: **for** $t = 1$ **to** T **do**
- 2: Get \mathbf{x}_t from an Online Convex Optimization algorithm
- 3: Receive stochastic subgradient \mathbf{g}_t such that $\mathbb{E}_t[\mathbf{g}_t] \in \partial F(\mathbf{x}_t)$
- 4: Pass loss $\ell_t(\mathbf{x}) = \langle \mathbf{g}_t, \mathbf{x} \rangle$ to Online Learning Algorithm
- 5: **end for**
- 6: **return** $\bar{\mathbf{x}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t$

Theorem

$$\mathbb{E}[F(\bar{\mathbf{x}}_T)] - F(\mathbf{u}) \leq \frac{\mathbb{E}[\text{Regret}_T(\mathbf{u})]}{T}, \quad \forall \mathbf{u} \in \mathcal{V}$$

Corollary: any result on regret translates to a result on convergence for stochastic optimization of convex functions

Proof.

$$\begin{aligned}
\mathbb{E}[F(\bar{\mathbf{x}}_T)] - F(\mathbf{u}) &\stackrel{\text{Jensen}}{\leq} \frac{1}{T} \sum_{t=1}^T (\mathbb{E}[F(\mathbf{x}_t)] - F(\mathbf{u})) \\
&\stackrel{\text{convexity}}{\leq} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\langle \mathbb{E}_t[\mathbf{g}_t], \mathbf{x}_t - \mathbf{u} \rangle] \\
&= \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\mathbb{E}_t[\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{u} \rangle]] \\
&\stackrel{\text{total expectation}}{=} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{u} \rangle] \\
&= \frac{\mathbb{E}[\text{Regret}_T(\mathbf{u})]}{T}
\end{aligned}$$

□

Example: Stochastic Subgradient Descent

Require: Feasible set $\mathcal{V} \subseteq \mathbb{R}^d$, $\mathbf{x}_1 \in \mathcal{V}$, $\eta = \frac{\alpha}{\sqrt{T}}$

1: **for** $t = 1$ **to** T **do**

2: Output $\mathbf{x}_t \in \mathcal{V}$

3: Receive stochastic subgradient \mathbf{g}_t such that $\mathbb{E}_t[\mathbf{g}_t] \in \partial F(\mathbf{x}_t)$

4: $\mathbf{x}_{t+1} = \Pi_{\mathcal{V}}(\mathbf{x}_t - \eta \mathbf{g}_t)$

5: **end for**

6: **return** $\bar{\mathbf{x}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t$

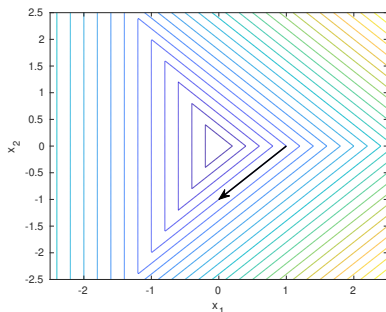
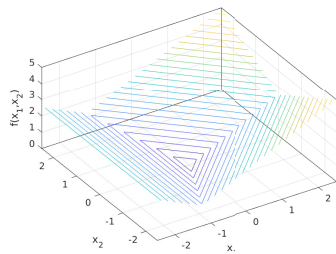
From the previous slides, we have

$$\mathbb{E}[F(\bar{\mathbf{x}}_T)] - F(\mathbf{x}^*) \leq \frac{1}{2\sqrt{T}} \left(\frac{\|\mathbf{x}_1 - \mathbf{x}^*\|_2^2}{\alpha} + \alpha \right)$$

In words, one pass of SGD minimizes the true risk

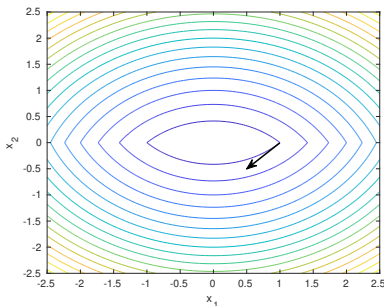
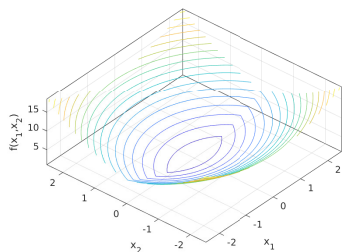
Beyond Online Subgradient Descent

Does Online Subgradient Descent Minimize the Functions? (1)



3D plot (left) and level sets (right) of $f(\mathbf{x}) = \max[-x_1, x_1 - x_2, x_1 + x_2]$. A negative subgradient is indicated by the black arrow

Does Online Subgradient Descent Minimize the Functions? (2)



3D plot (left) and level sets (right) of $f(\mathbf{x}) = \max[x_1^2 + (x_2 + 1)^2, x_1^2 + (x_2 - 1)^2]$. A negative subgradient is indicated by the black arrow

Understanding the Update of Online Subgradient Descent

$$\begin{aligned}\Pi_{\mathcal{V}}(\mathbf{x}_t - \eta_t \mathbf{g}_t) &= \underset{\mathbf{x} \in \mathcal{V}}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{x}_t + \eta_t \mathbf{g}_t\|_2^2 \\ &= \underset{\mathbf{x} \in \mathcal{V}}{\operatorname{argmin}} \|\eta_t \mathbf{g}_t\|_2^2 + 2\eta_t \langle \mathbf{g}_t, \mathbf{x} - \mathbf{x}_t \rangle + \|\mathbf{x}_t - \mathbf{x}\|_2^2 \\ &= \underset{\mathbf{x} \in \mathcal{V}}{\operatorname{argmin}} \underbrace{\ell_t(\mathbf{x}_t) + \langle \mathbf{g}_t, \mathbf{x} - \mathbf{x}_t \rangle}_{\text{Linear approximation of } \ell_t} + \frac{1}{2\eta_t} \underbrace{\|\mathbf{x}_t - \mathbf{x}\|_2^2}_{\text{Stay close to } \mathbf{x}_t}\end{aligned}$$

where $\Pi_{\mathcal{V}}$ is the Euclidean projection onto \mathcal{V} , i.e., $\Pi_{\mathcal{V}}(\mathbf{x}) = \underset{\mathbf{y} \in \mathcal{V}}{\operatorname{argmin}} \|\mathbf{x} - \mathbf{y}\|_2$

General Notion of Distances using Bregman Divergences

$$\operatorname{argmin}_{\mathbf{x} \in \mathcal{V}} \ell_t(\mathbf{x}_t) + \langle \mathbf{g}_t, \mathbf{x} - \mathbf{x}_t \rangle + \frac{1}{2\eta_t} \|\mathbf{x}_t - \mathbf{x}\|_2^2$$

Why the square Euclidean norm?

I can use general notion of distances, in particular **Bregman divergences**

Definition (Bregman Divergence [Bregman, 1967])

Let $\psi : X \rightarrow \mathbb{R}$ be strictly convex and differentiable on $\operatorname{int} X \neq \{\}$. The **Bregman Divergence** w.r.t. ψ is denoted by $B_\psi : X \times \operatorname{int} X \rightarrow \mathbb{R}$ defined as

$$B_\psi(\mathbf{x}; \mathbf{y}) = \psi(\mathbf{x}) - \psi(\mathbf{y}) - \langle \nabla \psi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle .$$

We start from the equivalent formulation of the OSD update

$$\mathbf{x}_{t+1} = \underset{\mathbf{x} \in \mathcal{V}}{\operatorname{argmin}} \ell_t(\mathbf{x}_t) + \langle \mathbf{g}_t, \mathbf{x} - \mathbf{x}_t \rangle + \frac{1}{2\eta_t} \|\mathbf{x}_t - \mathbf{x}\|_2^2$$

and we can change the last term with a Bregman Divergence

$$\mathbf{x}_{t+1} = \underset{\mathbf{x} \in \mathcal{V}}{\operatorname{argmin}} \ell_t(\mathbf{x}_t) + \langle \mathbf{g}_t, \mathbf{x} - \mathbf{x}_t \rangle + \frac{1}{\eta_t} B_\psi(\mathbf{x}; \mathbf{x}_t)$$

Require: $\psi : X \rightarrow \mathbb{R}$ strictly convex and differentiable on $\operatorname{int} X$, feasible set $\mathcal{V} \subseteq X \subseteq \mathbb{R}^d$, $\mathbf{x}_1 \in \operatorname{int} X \cap \mathcal{V}$

- 1: **for** $t = 1$ **to** T **do**
- 2: Output $\mathbf{x}_t \in \mathcal{V}$
- 3: Pay $\ell_t(\mathbf{x}_t)$
- 4: Set $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$
- 5: Set $\mathbf{x}_{t+1} \in \underset{\mathbf{x} \in \mathcal{V}}{\operatorname{argmin}} \langle \mathbf{g}_t, \mathbf{x} \rangle + \frac{1}{\eta_t} B_\psi(\mathbf{x}; \mathbf{x}_t)$
- 6: **end for**

Strongly Convex Functions

Definition

$f : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ is λ -strongly convex w.r.t. $\|\cdot\|$ if

$$f(\mathbf{x}) - f(\mathbf{y}) \leq \langle \mathbf{g}, \mathbf{x} - \mathbf{y} \rangle - \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{g} \in \partial f(\mathbf{x}).$$

Intuition: Lower bounded by a quadratic

Lemma (For OMD proof)

If ψ is λ -strongly convex w.r.t. $\|\cdot\|$ then $B_\psi(\mathbf{x}; \mathbf{y}) \geq \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|^2$

Lemma (For FTRL proof)

Let $f : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ closed, proper, subdifferentiable, and λ -strongly convex with respect to a norm $\|\cdot\|$ over its domain. Let $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x})$. Then, for all $\mathbf{x} \in \operatorname{dom} \partial f$, and $\mathbf{g} \in \partial f(\mathbf{x})$, we have

$$f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{1}{2\lambda} \|\mathbf{g}\|_*^2.$$

Theorem

Let ψ be λ -strongly convex w.r.t. $\|\cdot\|$. Pick any $\mathbf{x}_1 \in \text{int } X \cap \mathcal{V}$ and assume $\eta_t = \eta$, $t = 1, \dots, T$. Then, $\forall \mathbf{u} \in \mathcal{V}$, OMD satisfies

$$\sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) \leq \frac{B_\psi(\mathbf{u}; \mathbf{x}_1)}{\eta} + \frac{\eta}{2\lambda} \sum_{t=1}^T \|\mathbf{g}_t\|_*^2 - \frac{1}{\eta} B_\psi(\mathbf{u}; \mathbf{x}_{T+1}).$$

Proof.

One can show

$$\begin{aligned} \eta_t (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) &\leq \eta \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{u} \rangle \\ &\leq B_\psi(\mathbf{u}; \mathbf{x}_t) - B_\psi(\mathbf{u}; \mathbf{x}_{t+1}) - B_\psi(\mathbf{x}_{t+1}; \mathbf{x}_t) + \langle \eta_t \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}_{t+1} \rangle \end{aligned}$$

The last term can be bounded as

$$\langle \eta_t \mathbf{g}_t, \mathbf{x}_t - \mathbf{x}_{t+1} \rangle \leq \eta_t \|\mathbf{g}_t\|_* \|\mathbf{x}_t - \mathbf{x}_{t+1}\| \leq \frac{\eta_t^2 \|\mathbf{g}_t\|_*^2}{2\lambda} + \frac{\lambda}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2$$

From strong convexity of ψ , we get $-B_\psi(\mathbf{x}_{t+1}; \mathbf{x}_t) \leq -\frac{\lambda}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2$. Putting all together and summing over time, we get the stated bound. \square

Example: Online Subgradient Descent

- Set $\psi(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|_2^2$
- ψ is 1-strongly convex w.r.t. the L_2 norm
- Dual norm of L_2 is L_2

$$B(\mathbf{x}; \mathbf{y}) = \frac{1}{2}\|\mathbf{x}\|_2^2 - \frac{1}{2}\|\mathbf{y}\|_2^2 - \langle \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle = \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|_2^2$$

Regret for any \mathbf{u} :

$$\begin{aligned} \sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell(\mathbf{u})) &\leq \frac{B_\psi(\mathbf{u}; \mathbf{x}_1)}{\eta} + \frac{\eta}{2\lambda} \sum_{t=1}^T \|\mathbf{g}_t\|_*^2 \\ &= \frac{\|\mathbf{x}_1 - \mathbf{u}\|_2^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{g}_t\|_2^2 \end{aligned}$$

Example: Exponentiated Gradient (a.k.a. Hedge, EWA, etc.)

- Set $\mathcal{V} = \Delta^{d-1} := \{\mathbf{x} \in \mathbb{R}^d : x_i \geq 0, \|\mathbf{x}\|_1 = 1\}$
- Set $\psi(\mathbf{x}) = \sum_{i=1}^d x_i \ln x_i$
- ψ is 1-strongly convex w.r.t. the L_1 norm
- Dual norm of L_1 is L_∞
- Assume $\|\mathbf{g}_t\|_\infty \leq 1$

Require: $\eta > 0$

- 1: Set $\mathbf{x}_1 = [1/d, \dots, 1/d]$
- 2: **for** $t = 1$ **to** T **do**
- 3: Output $\mathbf{x}_t \in \Delta^{d-1}$
- 4: Pay $\ell_t(\mathbf{x}_t)$
- 5: Set $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$
- 6: $\mathbf{x}_{t+1, j} = \frac{x_{t, j} \exp(-\eta g_{t, j})}{\sum_{i=1}^d x_{t, i} \exp(-\eta g_{t, i})}$, $j = 1, \dots, d$
- 7: **end for**

Regret for any \mathbf{u} :
$$\sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell(\mathbf{u})) \leq \frac{B_\psi(\mathbf{u}; \mathbf{x}_1)}{\eta} + \frac{\eta}{2\lambda} \sum_{t=1}^T \|\mathbf{g}_t\|_*^2 \leq \frac{\ln d}{\eta} + \frac{\eta T}{2}$$

Set $\eta = \sqrt{\frac{2 \ln d}{T}}$ to obtain the upper bound of $\sqrt{2T \ln d}$

[Kivinen&Warmuth, 1997]

Follow-The-Regularized-Leader Algorithm

Require: Feasible set $\mathcal{V} \subseteq X \subseteq \mathbb{R}^d$, a sequence of regularizers

$$\psi_1, \dots, \psi_T : X \rightarrow \mathbb{R}$$

1: **for** $t = 1$ **to** T **do**

2: Output $\mathbf{x}_t \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{V}} \psi_t(\mathbf{x}) + \sum_{i=1}^{t-1} \ell_i(\mathbf{x})$

3: Receive $\ell_t : \mathcal{V} \rightarrow \mathbb{R}$ and pay $\ell_t(\mathbf{x}_t)$

4: **end for**

Lemma

Let $\psi_1, \dots, \psi_T : X \rightarrow \mathbb{R}$ be a sequence of regularization functions and $\mathcal{V} \subseteq X \subseteq \mathbb{R}^d$. Denote by $F_t(\mathbf{x}) = \psi_t(\mathbf{x}) + \sum_{i=1}^{t-1} \ell_i(\mathbf{x})$. Set $\mathbf{x}_t \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{V}} F_t(\mathbf{x})$. Then, for any $\mathbf{u} \in \mathbb{R}^d$, we have

$$\begin{aligned} \sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) &= \psi_{T+1}(\mathbf{u}) - \min_{\mathbf{x} \in \mathcal{V}} \psi_1(\mathbf{x}) + \sum_{t=1}^T [F_t(\mathbf{x}_t) - F_{t+1}(\mathbf{x}_{t+1}) + \ell_t(\mathbf{x}_t)] \\ &\quad + F_{T+1}(\mathbf{x}_{T+1}) - F_{T+1}(\mathbf{u}). \end{aligned}$$

Proof.

Just sum simplify the sums and use the fact that $F_1(\mathbf{x}_1) = \min_{\mathbf{x} \in \mathcal{V}} \psi_1(\mathbf{x})$. \square

An Explicit Regret with Strongly Convex Functions

Assume that $\psi_t + \sum_{i=1}^t \ell_i$ is λ_t strongly convex w.r.t. $\|\cdot\|$ and $\psi_{t+1}(\mathbf{x}) \geq \psi_t(\mathbf{x})$, we have

$$\begin{aligned} & \sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) \\ &= \psi_{T+1}(\mathbf{u}) - \min_{\mathbf{x} \in \mathcal{V}} \psi_1(\mathbf{x}) + \sum_{t=1}^T [F_t(\mathbf{x}_t) - F_{t+1}(\mathbf{x}_{t+1}) + \ell_t(\mathbf{x}_t)] + F_{T+1}(\mathbf{x}_{T+1}) - F_{T+1}(\mathbf{u}) \end{aligned}$$

Strong convexity

$$\leq \psi_{T+1}(\mathbf{u}) - \min_{\mathbf{x} \in \mathcal{V}} \psi_1(\mathbf{x}) + \sum_{t=1}^T \frac{\|\mathbf{g}_t\|_*^2}{2\lambda_t}$$

FTRL with Linearized Losses

- FTRL needs to solve a convex optimization problem at each step
- I can run FTRL with any sequence of losses
- I can also construct some losses
- For example, I might want to run FTRL on $\hat{\ell}_t(\mathbf{x}) = \ell_t(\mathbf{x}_t) + \langle \mathbf{g}_t, \mathbf{x} - \mathbf{x}_t \rangle$ where $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$

Require: A sequence of regularizers $\psi_1, \dots, \psi_T : X \rightarrow \mathbb{R}$

- 1: **for** $t = 1$ **to** T **do**
- 2: Output $\mathbf{x}_t \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{V}} \psi_t(\mathbf{x}) + \sum_{i=1}^{t-1} \langle \mathbf{g}_i, \mathbf{x} \rangle$
- 3: Pay $\ell_t(\mathbf{x}_t)$
- 4: Get $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$
- 5: **end for**

Same regret because

$$\sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) \leq \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{u} \rangle$$

- $\mathcal{V} = \mathbb{R}^d$
- $\psi_{t+1}(\mathbf{x}) = \frac{1}{2\eta_{t+1}} \|\mathbf{x}\|_2^2$
- $\mathbf{x}_{t+1} = \operatorname{argmin}_{\mathbf{x}} \frac{1}{2\eta_{t+1}} \|\mathbf{x}\|_2^2 + \sum_{i=1}^t \langle \mathbf{g}_i, \mathbf{x} \rangle = -\eta_{t+1} \sum_{i=1}^t \mathbf{g}_i$
- Compare it with OSD with $\mathbf{x}_1 = \mathbf{0}$: $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \mathbf{g}_t = -\sum_{i=1}^t \eta_i \mathbf{g}_i$

- **Important:** In FTRL the gradients are used with the same weight
- **Important:** In FTRL we don't take "jumps" of size η_t

Example: FTRL with Linearized Loss and Euclidean Regularization

- $\mathcal{V} = \mathbb{R}^d$
- $\psi(\mathbf{x}) = \frac{\gamma}{2} \|\mathbf{x}\|_2^2$
- ψ is γ -strongly convex w.r.t. L_2 norm
- Dual norm of L_2 norm is L_2 norm

$$\mathbf{x}_t = \underset{\mathbf{x} \in \mathcal{V}}{\operatorname{argmin}} \frac{\gamma}{2} \|\mathbf{x}\|_2^2 + \sum_{i=1}^{t-1} \langle \mathbf{g}_i, \mathbf{x} \rangle = \frac{-\sum_{i=1}^{t-1} \mathbf{g}_i}{\gamma}$$

$$\begin{aligned} \sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) &\leq \sum_{t=1}^T \langle \mathbf{g}_t, \mathbf{x}_t - \mathbf{u} \rangle \leq \psi_{T+1}(\mathbf{u}) - \min_{\mathbf{x} \in \mathcal{V}} \psi_1(\mathbf{x}) + \sum_{t=1}^T \frac{\|\mathbf{g}_t\|_*^2}{2\lambda_t} \\ &= \frac{\gamma}{2} \|\mathbf{u}\|_2^2 + \sum_{t=1}^T \frac{\|\mathbf{g}_t\|_2^2}{2\gamma} \end{aligned}$$

What is the optimal tuning of γ ?