

Online Prediction: Rademacher Averages via Burkholder's Method

Sasha Rakhlin

UPenn

Dec 16, 2017

Joint work with D. Foster and K. Sridharan

Outline

Motivation: Online Supervised Learning

Burkholder's Method

prediction paradigm

Data → Estimate Model → Make Prediction

prediction paradigm

Data → Estimate Model → Make Prediction

movies

users

					1		
	5						
			3				
					2		

movies

users

					1		
	5						
			3				
	?						
					2		

movies

users

					1		
	5						
			3				
	4						
					2		

movies

users

					1		
	5					?	
			3				
	4						
					2		

movies

users

					1		
	5					5	
			3				
	4						
					2		

Example 0: Bit Prediction

Suppose we want to predict a 0/1 sequence

$$y_1, y_2, \dots$$

Example 0: Bit Prediction

Suppose we want to predict a 0/1 sequence

$$y_1, y_2, \dots$$

If iid Bernoulli, then predicting majority $\mathbf{I}\{\bar{y}_t \geq .5\}$ ensures that proportion of correct predictions \bar{c}_t satisfies

$$\bar{c}_t \rightsquigarrow \max\{p, 1 - p\} \approx \max\{\bar{y}_t, 1 - \bar{y}_t\}$$

Example 0: Bit Prediction

Suppose we want to predict a 0/1 sequence

$$y_1, y_2, \dots$$

If iid Bernoulli, then predicting majority $\mathbf{I}\{\bar{y}_t \geq .5\}$ ensures that proportion of correct predictions \bar{c}_t satisfies

$$\bar{c}_t \rightsquigarrow \max\{p, 1 - p\} \approx \max\{\bar{y}_t, 1 - \bar{y}_t\}$$

More precisely:

$$\liminf_{n \rightarrow \infty} (\bar{c}_n - \max\{\bar{y}_n, 1 - \bar{y}_n\}) \geq 0 \quad \text{almost surely} \quad (*)$$

Example 0: Bit Prediction

Suppose we want to predict a 0/1 sequence

$$y_1, y_2, \dots$$

If iid Bernoulli, then predicting majority $\mathbf{I}\{\bar{y}_t \geq .5\}$ ensures that proportion of correct predictions \bar{c}_t satisfies

$$\bar{c}_t \rightsquigarrow \max\{p, 1-p\} \approx \max\{\bar{y}_t, 1-\bar{y}_t\}$$

More precisely:

$$\liminf_{n \rightarrow \infty} (\bar{c}_n - \max\{\bar{y}_n, 1-\bar{y}_n\}) \geq 0 \quad \text{almost surely} \quad (*)$$

Claim: *there is a method that ensures $(*)$ for an arbitrary sequence.*

Any idea how to do it? Majority will not work. Need randomized strategy.

Minimax vs Bayes prediction

By

D. Blackwell

University of California, Berkeley

Let $x = (x_1, x_2, \dots)$ be an infinite sequence of 0s and 1s, initially unknown to you. On day $n = 1, 2, \dots$, you observe $h_n = (x_1, \dots, x_{n-1})$, the first $n - 1$ terms of the sequence, and must predict x_n . What is a good prediction method, and how well can you do?

A *prediction method* p is just a function that associates with each finite sequence h of 0s and 1s a prediction $p(h) = 0$ or 1 , your prediction of the next x when you have observed history h . Denote by $w_n(p, x)$ the proportion of correct predictions that method p makes against sequence x in the first n days.

Looking for a p that is an x for which $w_n(p, x)$ approaches 50% as $n \rightarrow \infty$, with probability 1. A *random prediction method* is a function that chooses x_1 so

etc.

But we can improve

toss a fair coin every day, predicting 1 on heads, 0 on tails; the strong law of large numbers guarantees that, for every x , the proportion of correct predictions will approach 50% as $n \rightarrow \infty$, with probability 1. A *random prediction method* is a function that chooses x_1 so

The contrast between the minimax predictor p_0 and the Bayes predictor y is strong. The minimax predictor is not obvious, it is randomized, it satisfies (*) for every x , but the proof is not easy. The Bayes predictor is extremely obvious, it is not randomized, it satisfies (*) only for almost all x , and the proof is simple.

Online Supervised Learning

For $t = 1, \dots, n$
observe side info $x_t \in \mathcal{X}$
predict \hat{y}_t
observe outcome y_t

Online Supervised Learning

For $t = 1, \dots, n$
observe side info $x_t \in \mathcal{X}$
predict \hat{y}_t
observe outcome y_t

Goal:

$$\forall (x_t, y_t)_{t=1}^n, \quad \sum_{t=1}^n |\hat{y}_t - y_t| \leq$$

small if sequence is “nice”

Online Supervised Learning

For $t = 1, \dots, n$
observe side info $x_t \in \mathcal{X}$
predict \hat{y}_t
observe outcome y_t

Goal:

$$\forall (x_t, y_t)_{t=1}^n, \quad \sum_{t=1}^n |\hat{y}_t - y_t| \leq \phi(x_1, y_1, \dots, x_n, y_n)$$

Online Supervised Learning

For $t = 1, \dots, n$
observe side info $x_t \in \mathcal{X}$
predict \hat{y}_t
observe outcome y_t

Goal:

$$\forall (x_t, y_t)_{t=1}^n, \quad \sum_{t=1}^n |\hat{y}_t - y_t| \leq \inf_{w \in \mathcal{F}} \sum_{t=1}^n |\langle w, x_t \rangle - y_t| + C_n(x_1, \dots, x_n).$$

Online Supervised Learning

For $t = 1, \dots, n$
observe side info $x_t \in \mathcal{X}$
predict \hat{y}_t
observe outcome y_t

Goal:

$$\forall (x_t, y_t)_{t=1}^n, \quad \sum_{t=1}^n |\hat{y}_t - y_t| \leq \inf_{w \in \mathcal{F}} \sum_{t=1}^n |\langle w, x_t \rangle - y_t| + C_n(x_1, \dots, x_n).$$

If $C_n(x_1, \dots, x_n) = C_n$ is data-independent, we have full characterization.

Best possible: Empirical Rademacher

$$C_n(x_1, \dots, x_n) \sim \mathbb{E}_{\epsilon} \left\| \sum_{t=1}^n \epsilon_t x_t \right\|$$

Example: Matrix Completion.

Motivation: Online Supervised Learning

Proper learning: $\widehat{\mathbf{y}}_t = \langle \mathbf{w}_t, \mathbf{x}_t \rangle$,

$$\sum_{t=1}^n \underbrace{|\widehat{\mathbf{y}}_t - \mathbf{y}_t|}_{f_t(\mathbf{w}_t)} \leq \inf_{\mathbf{w} \in \mathcal{F}} \sum_{t=1}^n \underbrace{|\langle \mathbf{w}, \mathbf{x}_t \rangle - \mathbf{y}_t|}_{f_t(\mathbf{w})} + C_n(\mathbf{x}_1, \dots, \mathbf{x}_n)$$

Motivation: Online Supervised Learning

Proper learning: $\hat{\mathbf{y}}_t = \langle \mathbf{w}_t, \mathbf{x}_t \rangle$,

$$\sum_{t=1}^n \underbrace{|\hat{\mathbf{y}}_t - \mathbf{y}_t|}_{f_t(\mathbf{w}_t)} \leq \inf_{\mathbf{w} \in \mathcal{F}} \sum_{t=1}^n \underbrace{|\langle \mathbf{w}, \mathbf{x}_t \rangle - \mathbf{y}_t|}_{f_t(\mathbf{w})} + C_n(\mathbf{x}_1, \dots, \mathbf{x}_n)$$

Gradient Descent:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla f_t(\mathbf{w}_t)$$

Mirror Descent: different geometry

$$\nabla R(\mathbf{w}_{t+1}) = \nabla R(\mathbf{w}_t) - \eta \nabla f_t(\mathbf{w}_t)$$

(e.g. Exponential Weights Algorithm: $R(\mathbf{w}) = \sum \mathbf{w}(i) \log \mathbf{w}(i)$ is strongly convex w.r.t. ℓ_1 norm on simplex)

Motivation: Online Supervised Learning

Gradient/mirror descent with adaptive step size:

$$C_n \propto \sqrt{\sum_{t=1}^n \|\nabla f_t(\mathbf{w}_t)\|^2} = \sqrt{\sum_{t=1}^n \|\mathbf{x}_t\|^2}$$

Can be much worse than

$$\mathbb{E}_{\epsilon} \left\| \sum_{t=1}^n \epsilon_t \mathbf{x}_t \right\|.$$

Is this easy to fix?

Key issue: GD is not keeping the right *statistics* about the sequence

- ▶ Need additional information about geometry of functions: not just the size of gradients but also their “spread”
- ▶ Beyond usual notions of smoothness and strong convexity?

Key issue: GD is not keeping the right *statistics* about the sequence

- ▶ Need additional information about geometry of functions: not just the size of gradients but also their “spread”
- ▶ Beyond usual notions of smoothness and strong convexity?

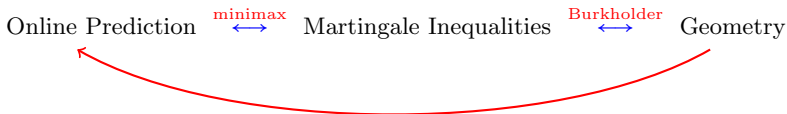
We would be searching in the dark if not for the connections:

Online Prediction $\overset{\text{minimax}}{\longleftrightarrow}$ Martingale Inequalities $\overset{\text{Burkholder}}{\longleftrightarrow}$ Geometry

Key issue: GD is not keeping the right *statistics* about the sequence

- ▶ Need additional information about geometry of functions: not just the size of gradients but also their “spread”
- ▶ Beyond usual notions of smoothness and strong convexity?

We would be searching in the dark if not for the connections:



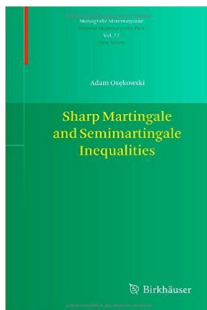
Outline

Motivation: Online Supervised Learning

Burkholder's Method

Reference:

Adam Osekowski, *Sharp Martingale and Semimartingale Inequalities*, 2012



Next: adaptation of some of these ideas to our setting.

$\epsilon_1, \dots, \epsilon_t, \dots$ i.i.d. Rademacher, $\mathcal{F}_t = \sigma(\epsilon_1, \dots, \epsilon_t)$.

X_1, \dots, X_t, \dots martingale difference sequence w.r.t. (\mathcal{F}_t) .

$$\mathbb{E}[X_t | \mathcal{F}_{t-1}] = 0$$

Note: X_t can be written as $\epsilon_t \mathbf{x}_t(\epsilon_{1:t-1})$ for some function \mathbf{x}_t .

A martingale inequality can be generically written as

$$\mathbb{E}B(X_1, \dots, X_n) \leq 0 \tag{1}$$

for some $B : \cup \mathcal{X}^n \rightarrow \mathbb{R}$.

A martingale inequality can be generically written as

$$\mathbb{E}B(X_1, \dots, X_n) \leq 0 \tag{1}$$

for some $B : \cup \mathcal{X}^n \rightarrow \mathbb{R}$.

(FRS'17+): (1) holds \forall martingale difference sequences (X_t) and all n

iff

there exists a function $U : \cup \mathcal{X}^n \rightarrow \mathbb{R}$ satisfying three properties:

A martingale inequality can be generically written as

$$\mathbb{E}B(X_1, \dots, X_n) \leq 0 \tag{1}$$

for some $B : \cup \mathcal{X}^n \rightarrow \mathbb{R}$.

(FRS'17+): (1) holds \forall martingale difference sequences (X_t) and all n

iff

there exists a function $U : \cup \mathcal{X}^n \rightarrow \mathbb{R}$ satisfying three properties:

1°. $U(x_1, \dots, x_t) \geq B(x_1, \dots, x_t)$ for all t

2°. $U(\cdot) \leq 0$

3°. for all x_1, \dots, x_t

$$\mathbb{E}_{\epsilon} U(x_1, \dots, x_{t-1}, \epsilon x_t) \leq U(x_1, \dots, x_{t-1})$$

proof

(\Leftarrow):

$$\mathbb{E}B(X_1, \dots, X_n) \stackrel{1^\circ}{\leq} \mathbb{E}U(X_1, \dots, X_n) \stackrel{3^\circ}{\leq} \mathbb{E}U(X_1, \dots, X_{n-1}) \dots \leq U(\cdot) \stackrel{2^\circ}{\leq} 0.$$

proof

(\Leftarrow) :

$$\mathbb{E}B(X_1, \dots, X_n) \stackrel{1^\circ}{\leq} \mathbb{E}U(X_1, \dots, X_n) \stackrel{3^\circ}{\leq} \mathbb{E}U(X_1, \dots, X_{n-1}) \dots \leq U(\cdot) \stackrel{2^\circ}{\leq} 0.$$

(\Rightarrow) : Define

$$U^*(x_1, \dots, x_t) \triangleq \sup_{n \geq t, (X)_{t+1}^n} \mathbb{E}B(x_1, \dots, x_t, X_{t+1}, \dots, X_n).$$

Claim: U^* satisfies $1^\circ - 3^\circ$.

Claim: u^* is the smallest function that satisfies 1° – 3°.

Proof:

For any $n \geq t$, $(X)_{t+1}^n$,

$$\mathbb{E}B(x_1, \dots, x_t, X_{t+1}, \dots, X_n) \leq \mathbb{E}u'(x_1, \dots, x_t, X_{t+1}, \dots, X_n) \leq u'(x_1, \dots, x_t)$$

Hence, $u^*(x_1, \dots, x_t) \leq u'(x_1, \dots, x_t)$.

Example 1: Smoothness/Strong Convexity

Let $\|\cdot\|$ be some norm on \mathcal{X} . Martingale inequality corresponding to smoothness:

$$\mathbb{E} \left\| \sum_{t=1}^n \mathbf{x}_t \right\|^2 \leq K \cdot \mathbb{E} \sum_{t=1}^n \|\mathbf{x}_t\|^2$$

Hence,

$$B(\mathbf{x}_1, \dots, \mathbf{x}_n) = B\left(\sum_{t=1}^n \mathbf{x}_t, \sum_{t=1}^n \|\mathbf{x}_t\|^2\right) = \left\| \sum_{t=1}^n \mathbf{x}_t \right\|^2 - K \sum_{t=1}^n \|\mathbf{x}_t\|^2$$

Example 1: Smoothness/Strong Convexity

Let $\|\cdot\|$ be some norm on \mathcal{X} . Martingale inequality corresponding to smoothness:

$$\mathbb{E} \left\| \sum_{t=1}^n \mathbf{x}_t \right\|^2 \leq K \cdot \mathbb{E} \sum_{t=1}^n \|\mathbf{x}_t\|^2$$

Hence,

$$B(\mathbf{x}_1, \dots, \mathbf{x}_n) = B\left(\sum_{t=1}^n \mathbf{x}_t, \sum_{t=1}^n \|\mathbf{x}_t\|^2\right) = \left\| \sum_{t=1}^n \mathbf{x}_t \right\|^2 - K \sum_{t=1}^n \|\mathbf{x}_t\|^2$$

Optimal \mathbf{u}^* inherits

$$U^*(\mathbf{x}, \alpha^2) = U^*(\mathbf{x}, 0) - K\alpha^2.$$

Example 1: Smoothness/Strong Convexity

Let $\|\cdot\|$ be some norm on \mathcal{X} . Martingale inequality corresponding to smoothness:

$$\mathbb{E} \left\| \sum_{t=1}^n \mathbf{x}_t \right\|^2 \leq K \cdot \mathbb{E} \sum_{t=1}^n \|\mathbf{x}_t\|^2$$

Hence,

$$B(\mathbf{x}_1, \dots, \mathbf{x}_n) = B\left(\sum_{t=1}^n \mathbf{x}_t, \sum_{t=1}^n \|\mathbf{x}_t\|^2\right) = \left\| \sum_{t=1}^n \mathbf{x}_t \right\|^2 - K \sum_{t=1}^n \|\mathbf{x}_t\|^2$$

Optimal \mathbf{u}^* inherits

$$U^*(\mathbf{x}, \alpha^2) = U^*(\mathbf{x}, 0) - K\alpha^2.$$

Restricted concavity 3° is

$$\mathbb{E}_{\epsilon} U^*(\mathbf{x} + \epsilon \mathbf{y}, \alpha^2 + \|\mathbf{y}\|^2) \leq U^*(\mathbf{x}, \alpha^2).$$

Example 1: Smoothness/Strong Convexity

Let $\|\cdot\|$ be some norm on \mathcal{X} . Martingale inequality corresponding to smoothness:

$$\mathbb{E} \left\| \sum_{t=1}^n \mathbf{x}_t \right\|^2 \leq K \cdot \mathbb{E} \sum_{t=1}^n \|\mathbf{x}_t\|^2$$

Hence,

$$B(\mathbf{x}_1, \dots, \mathbf{x}_n) = B\left(\sum_{t=1}^n \mathbf{x}_t, \sum_{t=1}^n \|\mathbf{x}_t\|^2\right) = \left\| \sum_{t=1}^n \mathbf{x}_t \right\|^2 - K \sum_{t=1}^n \|\mathbf{x}_t\|^2$$

Optimal \mathbf{u}^* inherits

$$U^*(\mathbf{x}, \alpha^2) = U^*(\mathbf{x}, 0) - K\alpha^2.$$

Restricted concavity 3° is

$$\mathbb{E}_{\epsilon} U^*(\mathbf{x} + \epsilon \mathbf{y}, \alpha^2 + \|\mathbf{y}\|^2) \leq U^*(\mathbf{x}, \alpha^2).$$

Corollary: $\mathbf{x} \mapsto U^*(\mathbf{x}, 0)$ is smooth wrt $\|\cdot\|$ (and its dual is strongly cvx).

Proof:

$$\begin{aligned}\Phi(x) &= U^*(x, 0) \\ &\geq \mathbb{E}_{\epsilon} U^*(x + \epsilon y, \|y\|^2) \\ &= \frac{1}{2}(U^*(x + y, 0) - K \|y\|^2) + \frac{1}{2}(U^*(x - y, 0) - K \|y\|^2) \\ &= \frac{1}{2}\Phi(x + y) + \frac{1}{2}\Phi(x - y) - K \|y\|^2\end{aligned}$$

Example 2

Hilbert space:

$$\mathbb{E} \left\| \sum_{t=1}^n X_t \right\| \leq 2\mathbb{E} \sqrt{\sum_{t=1}^n \|X_t\|^2}$$

Hence,

$$B(x, \alpha) = \|x\| - 2y$$

and interested in

$$B\left(\sum_{t=1}^n x_t, \sqrt{\sum_{t=1}^n \|x_t\|^2}\right)$$

Example 2

Hilbert space:

$$\mathbb{E} \left\| \sum_{t=1}^n X_t \right\| \leq 2 \mathbb{E} \sqrt{\sum_{t=1}^n \|X_t\|^2}$$

Hence,

$$B(x, a) = \|x\| - 2a$$

and interested in

$$B\left(\sum_{t=1}^n x_t, \sqrt{\sum_{t=1}^n \|x_t\|^2}\right)$$

Elementary algebra gives

$$U(x, a) = \begin{cases} -\sqrt{2a^2 - \|x\|^2}, & a \geq \|x\| \\ \|x\| - 2a, & a < \|x\| \end{cases}$$

satisfies concavity property

$$U(x + d, \sqrt{a^2 + d^2}) \leq U(x, a) + U_x(x, a)d$$

Example 2

Hilbert space:

$$\mathbb{E} \left\| \sum_{t=1}^n X_t \right\| \leq 2 \mathbb{E} \sqrt{\sum_{t=1}^n \|X_t\|^2}$$

Hence,

$$B(x, a) = \|x\| - 2a$$

and interested in

$$B\left(\sum_{t=1}^n x_t, \sqrt{\sum_{t=1}^n \|x_t\|^2}\right)$$

Elementary algebra gives

$$U(x, a) = \begin{cases} -\sqrt{2a^2 - \|x\|^2}, & a \geq \|x\| \\ \|x\| - 2a, & a < \|x\| \end{cases}$$

satisfies concavity property

$$U(x+d, \sqrt{a^2+d^2}) \leq U(x, a) + \frac{x \cdot d}{\sqrt{2a^2 - \|x\|^2}}$$

Example 3: Empirical Rademacher

(FRS'16): Empirical Rademacher bound for Online Supervised Learning is possible if (and close to “iff”)

$$\mathbb{E}_{\epsilon} \left\| \sum_{t=1}^n \epsilon_t \mathbf{x}_t \right\|^p \leq K \mathbb{E}_{\epsilon, \epsilon'} \left\| \sum_{t=1}^n \epsilon'_t \mathbf{x}_t \right\|^p$$

where \mathbf{x}_t is \mathcal{F}_{t-1} -measurable, $p \geq 1$. *Decoupling inequality.*

Example 3: Empirical Rademacher

(FRS'16): Empirical Rademacher bound for Online Supervised Learning is possible if (and close to “iff”)

$$\mathbb{E}_{\epsilon} \left\| \sum_{t=1}^n \epsilon_t \mathbf{x}_t \right\|^p \leq K \mathbb{E}_{\epsilon, \epsilon'} \left\| \sum_{t=1}^n \epsilon'_t \mathbf{x}_t \right\|^p$$

where \mathbf{x}_t is \mathcal{F}_{t-1} -measurable, $p \geq 1$. *Decoupling inequality*.

Two-sided version of above is equivalent to *deterministic UMD*

$$\forall \sigma_1, \dots, \sigma_n \in \{\pm 1\}, \quad \mathbb{E} \left\| \sum_{t=1}^n \sigma_t X_t \right\|^p \leq K \mathbb{E} \left\| \sum_{t=1}^n X_t \right\|^p$$

Example 3: Empirical Rademacher

(FRS'16): Empirical Rademacher bound for Online Supervised Learning is possible if (and close to “iff”)

$$\mathbb{E}_{\epsilon} \left\| \sum_{t=1}^n \epsilon_t \mathbf{x}_t \right\|^p \leq K \mathbb{E}_{\epsilon, \epsilon'} \left\| \sum_{t=1}^n \epsilon'_t \mathbf{x}_t \right\|^p$$

where \mathbf{x}_t is \mathcal{F}_{t-1} -measurable, $p \geq 1$. *Decoupling inequality*.

Two-sided version of above is equivalent to *deterministic UMD*

$$\forall \sigma_1, \dots, \sigma_n \in \{\pm 1\}, \quad \mathbb{E} \left\| \sum_{t=1}^n \sigma_t X_t \right\|^p \leq K \mathbb{E} \left\| \sum_{t=1}^n X_t \right\|^p$$

This gives

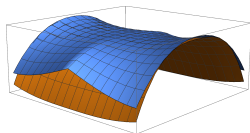
$$B(x_1, \dots, x_n) = B\left(\sum_{t=1}^n x_t, \sum_{t=1}^n \sigma_t x_t\right) = \left\| \sum_{t=1}^n \sigma_t x_t \right\|^p - K \left\| \sum_{t=1}^n x_t \right\|^p$$

Example 3: Empirical Rademacher

Then \mathcal{U} satisfying $1^\circ - 3^\circ$ is called Burkholder function. Property 3° reads

$$\mathbb{E}_\epsilon \mathcal{U}(x + \epsilon z, y + \epsilon z) \leq \mathcal{U}(x, y)$$

which is equivalent to *zigzag concavity*.



We can now go back and use \mathcal{U} to derive an algorithm for Online Supervised Learning with Empirical Rademacher regret bound.

Back to Online Supervised Learning

Enough to solve linearized problem

$$\sum_{t=1}^n \ell'_t \cdot \widehat{\mathbf{y}}_t \leq \min_{\mathbf{w} \in \mathcal{F}} \sum_{t=1}^n \ell'_t \cdot \langle \mathbf{w}, \mathbf{x}_t \rangle + C_n(\mathbf{x}_1, \dots, \mathbf{x}_n)$$

Back to Online Supervised Learning

Enough to solve linearized problem

$$\sum_{t=1}^n \ell'_t \cdot \widehat{\mathbf{y}}_t \leq \min_{\mathbf{w} \in \mathcal{F}} \sum_{t=1}^n \ell'_t \cdot \langle \mathbf{w}, \mathbf{x}_t \rangle + C_n(\mathbf{x}_1, \dots, \mathbf{x}_n)$$

which can be written as

$$\sum_{t=1}^n \ell'_t \cdot \widehat{\mathbf{y}}_t + \underbrace{\left\| \sum_{t=1}^n \ell'_t \mathbf{x}_t \right\| - C \mathbb{E}_{\epsilon} \left\| \sum_{t=1}^n \epsilon_t \ell'_t \mathbf{x}_t \right\|}_{\leq \mathbb{E}_{\epsilon} \mathcal{U}(\sum_{t=1}^n \ell'_t \mathbf{x}_t, \sum_{t=1}^n \epsilon_t \ell'_t \mathbf{x}_t)} \leq 0$$

Last step:

$$\min_{\widehat{\mathbf{y}}_n} \max_{\ell'_n \in [-1,1]} \left\{ \ell'_n \cdot \widehat{\mathbf{y}}_n + \mathbb{E}_{\epsilon} \mathcal{U} \left(\sum_{t=1}^n \ell'_t x_t, \sum_{t=1}^n \epsilon_t \ell'_t x_t \right) \right\}$$

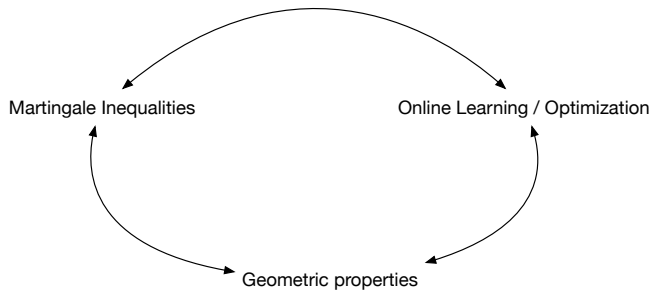
Choosing $\widehat{\mathbf{y}}_n = -G'(0)$ for

$$G(\alpha) = \mathbb{E}_{\sigma} \mathcal{U} \left(\sum_{t=1}^{n-1} \ell'_t x_t + \alpha x_t, \sum_{t=1}^{n-1} \epsilon_t \ell'_t x_t + \sigma \alpha x_t \right)$$

ensures

$$-\ell'_n \cdot G'(0) + G(\ell_n) \leq G(0)$$

by diagonal concavity and yields clean recursion.



Conclusions

- ▶ Gradient/Mirror Descent does not keep the right “statistics” about the sequence.
- ▶ Strong convexity/smoothness is not enough as a geometric primitive.
- ▶ Can find the right primitive by exploiting connections between probabilistic inequalities and geometry.