

# ONLINE ACCELERATION OF EXPONENTIAL WEIGHTS

---

Pierre Gaillard

February 10, 2017

INRIA Paris – SIERRA

This is joint work with Olivier Wintenberger

SETTING

---

Step by step minimization of i.i.d. convex loss functions<sup>1</sup>  $\ell_1, \dots, \ell_n : \mathbb{R}^d \rightarrow \mathbb{R}$ .

### Assumption 1 (strongly convex risk)

$$\exists \alpha > 0, \theta^* \in \mathbb{R}^d, \forall \theta \in \mathbb{R}^d \quad \alpha \|\theta - \theta^*\|_2^2 \leq \mathbb{E}[\ell_t(\theta) - \ell_t(\theta^*)].$$

### Setting: for each $t = 1, \dots, n$

- the learner provides  $\hat{\theta}_{t-1} \in \mathbb{R}^d$  based on past gradients  $\nabla \ell_s(\hat{\theta}_{s-1})$  for  $s \leq t-1$
- the environment reveals  $\nabla \ell_t(\hat{\theta}_{t-1})$

**Goal:** minimize the average risk:

$$\text{Risk}_n(\hat{\theta}_{0:(n-1)}) = \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\ell_t](\hat{\theta}_{t-1}) - \mathbb{E}[\ell_t](\theta^*).$$

### Remark 1

- non-Lipschitz gradients
- only the risk needs to be strongly convex (pinball loss)

<sup>1</sup> Cesa-Bianchi and Lugosi, *Prediction, Learning, and Games*, 2006.

Step by step minimization of i.i.d. convex loss functions<sup>1</sup>  $\ell_1, \dots, \ell_n : \mathbb{R}^d \rightarrow \mathbb{R}$ .

### Assumption 1 (strongly convex risk)

$$\exists \alpha > 0, \theta^* \in \mathbb{R}^d, \forall \theta \in \mathbb{R}^d \quad \alpha \|\theta - \theta^*\|_2^2 \leq \mathbb{E}[\ell_t(\theta) - \ell_t(\theta^*)].$$

**Setting:** for each  $t = 1, \dots, n$

- the learner provides  $\hat{\theta}_{t-1} \in \mathbb{R}^d$  based on past gradients  $\nabla \ell_s(\hat{\theta}_{s-1})$  for  $s \leq t-1$
- the environment reveals  $\nabla \ell_t(\hat{\theta}_{t-1})$

**Goal:** minimize the average risk:

$$\text{Risk}_n(\hat{\theta}_{0:(n-1)}) = \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\ell_t](\hat{\theta}_{t-1}) - \mathbb{E}[\ell_t](\theta^*).$$

**Remark 2** By convexity of the risk, the averaging  $\bar{\theta}_{n-1} := (1/n) \sum_{t=1}^n \hat{\theta}_{t-1}$  has an instantaneous risk upper-bounded by the cumulative risk

$$\text{Risk}(\bar{\theta}_{n-1}) := \mathbb{E}[\ell_n](\bar{\theta}_{n-1}) - \mathbb{E}[\ell_n](\theta^*) \stackrel{\text{Jensen}}{\leq} \text{Risk}_n(\hat{\theta}_{0:(n-1)}).$$

<sup>1</sup> Cesa-Bianchi and Lugosi, *Prediction, Learning, and Games*, 2006.

Step by step minimization of i.i.d. convex loss functions<sup>1</sup>  $\ell_1, \dots, \ell_n : \mathbb{R}^d \rightarrow \mathbb{R}$ .

### Assumption 1 (strongly convex risk)

$$\exists \alpha > 0, \theta^* \in \mathbb{R}^d, \forall \theta \in \mathbb{R}^d \quad \alpha \|\theta - \theta^*\|_2^2 \leq \mathbb{E}[\ell_t(\theta) - \ell_t(\theta^*)].$$

$\theta^*$  is  $d_0$ -sparse

$$\|\theta^*\|_1 \leq U$$

**Setting:** for each  $t = 1, \dots, n$

- the learner provides  $\hat{\theta}_{t-1} \in \mathbb{R}^d$  based on past gradients  $\nabla \ell_s(\hat{\theta}_{s-1})$  for  $s \leq t-1$
- the environment reveals  $\nabla \ell_t(\hat{\theta}_{t-1})$

**Goal:** minimize the average risk:

$$\text{Risk}_n(\hat{\theta}_{0:(n-1)}) = \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\ell_t(\hat{\theta}_{t-1}) - \mathbb{E}[\ell_t](\theta^*)].$$

**Remark 2** By convexity of the risk, the averaging  $\bar{\theta}_{n-1} := (1/n) \sum_{t=1}^n \hat{\theta}_{t-1}$  has an instantaneous risk upper-bounded by the cumulative risk

$$\text{Risk}(\bar{\theta}_{n-1}) := \mathbb{E}[\ell_n](\bar{\theta}_{n-1}) - \mathbb{E}[\ell_n](\theta^*) \stackrel{\text{Jensen}}{\leq} \text{Risk}_n(\hat{\theta}_{0:(n-1)}).$$

<sup>1</sup> Cesa-Bianchi and Lugosi, *Prediction, Learning, and Games*, 2006.

Step by step minimization of i.i.d. convex loss functions  $\ell_1, \dots, \ell_n : \mathbb{R}^d \rightarrow \mathbb{R}$ .

## Assumption 1 (strongly convex risk)

$$\exists \alpha > 0, \theta^* \in \mathbb{R}^d, \forall \theta \in \mathbb{R}^d \quad \alpha \|\theta - \theta^*\|_2^2 \leq \mathbb{E}[\ell_t(\theta) - \ell_t(\theta^*)].$$

$\theta^*$  is  $d_0$ -sparse

$$\|\theta^*\|_1 \leq U$$

**Setting:** for each  $t = 1, \dots, n$

- the learner provides  $\hat{\theta}_{t-1} \in \mathbb{R}^d$  based on past gradients  $\nabla \ell_s(\hat{\theta}_{s-1})$  for  $s \leq t-1$
- the environment reveals  $\nabla \ell_t(\hat{\theta}_{t-1})$

**Goal:** minimize the average risk:

$$\text{Risk}_n(\hat{\theta}_{0:(n-1)}) = \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\ell_t](\hat{\theta}_{t-1}) - \mathbb{E}[\ell_t](\theta^*).$$

**Result**

$$\text{Risk}_n(\hat{\theta}_{0:(n-1)}) \lesssim \min \left\{ \frac{B^2 d_0 \log d \log n}{\alpha n}, UB \sqrt{\frac{\log d}{n}} \right\}$$

where  $B \geq \max_{\|\theta\|_1 \leq 2U} \|\nabla \ell_t(\theta)\|_\infty$ .

Step by step minimization of i.i.d. convex loss functions  $\ell_1, \dots, \ell_n : \mathbb{R}^d \rightarrow \mathbb{R}$ .

## Assumption 1 (strongly convex risk)

$$\exists \alpha > 0, \theta^* \in \mathbb{R}^d, \forall \theta \in \mathbb{R}^d \quad \alpha \|\theta - \theta^*\|_2^2 \leq \mathbb{E}[\ell_t(\theta) - \ell_t(\theta^*)].$$

$\theta^*$  is  $d_0$ -sparse

$$\|\theta^*\|_1 \leq U$$

**Setting:** for each  $t = 1, \dots, n$

- the learner provides  $\hat{\theta}_{t-1} \in \mathbb{R}^d$  based on past gradients  $\nabla \ell_s(\hat{\theta}_{s-1})$  for  $s \leq t-1$
- the environment reveals  $\nabla \ell_t(\hat{\theta}_{t-1})$

**Goal:** minimize the average risk:

$$\text{Risk}_n(\hat{\theta}_{0:(n-1)}) = \frac{1}{n} \sum_{t=1}^n \mathbb{E}[\ell_t](\hat{\theta}_{t-1}) - \mathbb{E}[\ell_t](\theta^*).$$

**Result**

$$\text{Risk}_n(\hat{\theta}_{0:(n-1)}) \lesssim \min \left\{ \frac{B^2 d_0 \log d \log n}{\alpha n}, UB \sqrt{\frac{\log d}{n}} \right\}$$

where  $B \geq \max_{\|\theta\|_1 \leq 2U} \|\nabla \ell_t(\theta)\|_\infty$ .

Fast rate : better for large  $n, \alpha$

Slow rate : better for small  $n, \alpha$

Procedure	Sequential	Rate	Polynomial
Lasso <sup>1</sup>	✗	$\frac{d_0 \log d}{\alpha n}$	✓
EWA + sparsity pattern <sup>2</sup>	✗	$\frac{d_0 \log d}{n}$	✗
SeqSEW <sup>3</sup>	✓	$\frac{d_0 \log d}{n}$	✗
$\ell_1$ -RDA method <sup>4</sup>	✓	$\frac{d}{n}$	✓
SAEW	✓	$\frac{d_0 \log d}{\alpha n}$	✓

<sup>1</sup> Bunea, Tsybakov, and Wegkamp, "Aggregation for Gaussian regression", 2007.

<sup>2</sup> Rigollet and Tsybakov, "Exponential screening and optimal rates of sparse estimation", 2011.

<sup>3</sup> Gerchinovitz, "Sparsity regret bounds for individual sequences in online linear regression", 2013.

<sup>4</sup> Xiao, "Dual averaging methods for regularized stochastic learning and online optimization", 2010.

CONVEX OPTIMIZATION IN THE  $\ell_1$ -BALL  
WITH SLOW RATE

---

**Goal:** perform online optimization in

$$\mathcal{B}_1(\theta_{\text{center}}, \varepsilon) := \{\theta \in \mathbb{R}^d : \|\theta - \theta_{\text{center}}\|_1 \leq \varepsilon\}$$

We define the  $2d$  corners of the  $\ell_1$ -ball  $e_k = \theta_{\text{center}} \pm \varepsilon(0, \dots, 0, 1, 0, \dots, 0)$

### The exponentially gradient forecaster ( $EG^\pm$ )<sup>5</sup>

At each forecasting instance  $t \geq 1$ ,

- assign to each corner  $e_k$  the weight

$$\hat{p}_{k,t-1} = \frac{\exp\left(-\eta \sum_{s=1}^{t-1} \nabla \ell_s(\hat{\theta}_{s-1})^\top e_k\right)}{\sum_{j=1}^{2d} \exp\left(-\eta \sum_{s=1}^{t-1} \nabla \ell_s(\hat{\theta}_{s-1})^\top e_j\right)}$$

- form parameter  $\hat{\theta}_{t-1} = \sum_{k=1}^{2d} \hat{p}_{k,t-1} e_k$

**Performance:** bound on the average regret, if  $\theta^* \in \mathcal{B}_1(\theta_{\text{center}}, \varepsilon)$  for  $\eta$  well-tuned

$$\sum_{t=1}^n \ell_t(\hat{\theta}_{t-1}) - \ell_t(\theta^*) \lesssim \varepsilon B \sqrt{\frac{\log d}{n}}.$$

The learning rate  $\eta$  can be tuned online (doubling trick,  $\eta_t$ ).

<sup>5</sup> Kivinen and Warmuth, "Exponentiated Gradient Versus Gradient Descent for Linear Predictors", 1997.

### Lemma (Hoeffding)

If  $X$  is a random variable with  $|X| \leq B$ . Then,

$$\forall \eta \in \mathbb{R}, \quad \mathbb{E}[X] \leq -\frac{1}{\eta} \log \left( \mathbb{E}[e^{-\eta X}] \right) + \frac{\eta B}{4}.$$

1. Upper bound the instantaneous gradient

$$\begin{aligned} \nabla \ell_t(\hat{\theta}_{t-1})^\top \hat{\theta}_{t-1} &\stackrel{\text{def. of } \hat{\theta}_{t-1}}{=} \sum_{k=1}^{2d} \hat{\rho}_{k,t-1} (\nabla \ell_t(\hat{\theta}_{t-1})^\top e_k) \\ &\stackrel{\text{Hoeffding}}{\leq} -\frac{1}{\eta} \log \left( \sum_{k=1}^K \hat{\rho}_{k,t} e^{-\eta \nabla \ell_t(\hat{\theta}_{t-1})^\top e_k} \right) + \frac{\eta B}{4} \\ &\stackrel{\text{def. of } \hat{\rho}_{k,t+1}}{=} -\frac{1}{\eta} \log \left( \frac{\hat{\rho}_{k,t}}{\hat{\rho}_{k,t+1}} e^{-\eta \nabla \ell_t(\hat{\theta}_{t-1})^\top e_k} \right) + \frac{\eta B}{4} \\ &= \nabla \ell_t(\hat{\theta}_{t-1})^\top e_k + \frac{1}{\eta} \log \frac{\hat{\rho}_{k,t+1}}{\hat{\rho}_{k,t}} + \frac{\eta B}{4}. \end{aligned}$$

2. Sum over all  $t$ , the sum telescopes

$$\begin{aligned} \sum_{t=1}^n \ell_t(\hat{\theta}_{t-1}) - \ell_t(\hat{\theta}^*) &\stackrel{\text{Jensen}}{\leq} \sum_{t=1}^n \nabla \ell_t(\hat{\theta}_{t-1})^\top (\hat{\theta}_{t-1} - \hat{\theta}^*) \leq \max_{k=1, \dots, 2d} \left\{ \sum_{t=1}^n \nabla \ell_t(\hat{\theta}_{t-1})^\top (\hat{\theta}_{t-1} - e_k) \right\} \\ &\leq \max_{k=1, \dots, 2d} \left\{ \frac{1}{\eta} \log \frac{\hat{\rho}_{k,n}}{\hat{\rho}_{k,0}} + \frac{\eta B n}{4} \right\} \leq \frac{\log(2d)}{\eta} + \frac{\eta B n}{4} \end{aligned}$$

## Theorem

Let  $0 < \delta < 1$ , then if  $\theta^* \in \mathcal{B}_1(\theta^*, \varepsilon)$ ,  $EG^\pm$  satisfies

$$\alpha \|\bar{\theta}_{n-1} - \theta^*\|_2^2 \stackrel{\text{strong convexity}}{\leq} \mathbb{E}[l_n](\bar{\theta}_{n-1}) - \mathbb{E}[l_n](\theta^*) \lesssim \frac{\varepsilon B \sqrt{\log(d/\delta)}}{\sqrt{n}}$$

where  $\bar{\theta}_{n-1} = \frac{1}{n} \sum_{t=1}^n \hat{\theta}_{t-1}$

We observe the slow rate on the risk of order  $UB\sqrt{\log(d)/n}$ .

**Proof:**

- Hoeffding inequality for martingal : with high probability

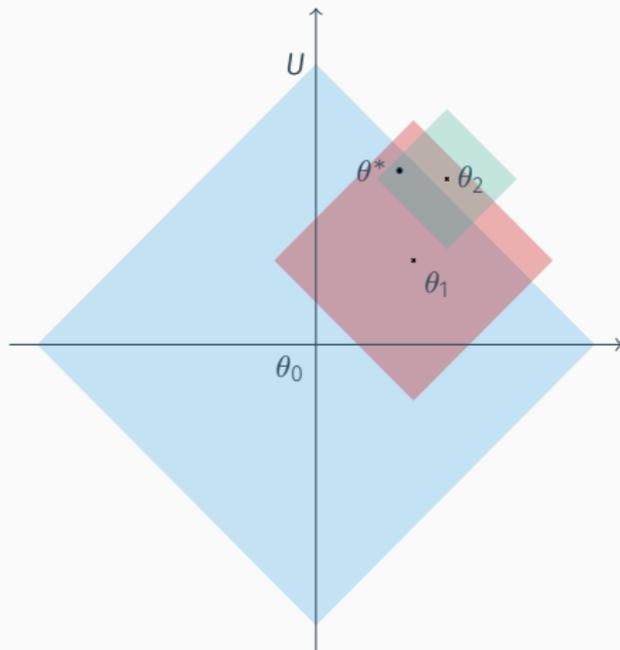
$$\sum_{t=1}^n \mathbb{E}[l_t](\hat{\theta}_{t-1}) - \mathbb{E}[l_t](\theta^*) \lesssim \sum_{t=1}^n \ell_t(\hat{\theta}_{t-1}) - \mathbb{E}[l_t](\theta^*) + \sqrt{n \log(1/\delta)}$$

- Jensen's inequality (convex i.i.d losses):

$$\mathbb{E}[l_n](\bar{\theta}_{n-1}) - \mathbb{E}[l_n](\theta^*) \leq \frac{1}{n} \sum_{t=1}^n \mathbb{E}[l_t](\hat{\theta}_{t-1}) - \mathbb{E}[l_t](\theta^*)$$

ACCELERATION : FROM SLOW RATE TO  
FAST RATE

---



# ACCELERATION: FROM SLOW RATE $1/\sqrt{n}$ TO FAST RATE $1/n$

## Algorithm

Parameters:  $U, B, \alpha, \delta > 0$

Initialization:  $\theta_{\text{center}} = 0 \in \mathbb{R}^d, t_1 = 1$

For sessions  $i \geq 1$ ,

- Start a new  $EG^\pm$  in  $\mathcal{B}_1(\theta_{\text{center}}, U2^{-i})$  for  $t \geq t_i$
- Get the high probability  $\ell_1$ -ball for  $\theta^*$

$$\|\bar{\theta}_{t-1} - \theta^*\|_1^2 \leq d \|\bar{\theta}_{t-1} - \theta^*\|_2^2 \lesssim \frac{dU2^{-i}B\sqrt{\log(d/\delta)}}{\alpha\sqrt{t-t_i}} =: C(U, t)^2$$

- Define  $t_{i+1} \geq t_i$  as the first time such that  $C(U, t_{i+1}) \leq U2^{-(i+1)}$ , i.e., for

$$\sqrt{t_{i+1} - t_i} \approx \frac{4dB\sqrt{\log(d/\delta)}}{U2^{-i}\alpha}$$

- $\theta_{\text{center}} \leftarrow \bar{\theta}_{t_{i+1}-1}$

After  $n$  time steps, we will have the slow rate high probability bound

$$\mathbb{E}[l_n](\bar{\theta}_{n-1}) - \mathbb{E}[l_n](\hat{\theta}_{t-1}) \lesssim \frac{U2^{-i}B\sqrt{\log(d/\delta)}}{\sqrt{n}}$$

but with  $U2^{-i} \approx \frac{dB\sqrt{\log(d/\delta)}}{\alpha\sqrt{n}}$ .

# ACCELERATION: FROM SLOW RATE $1/\sqrt{n}$ TO FAST RATE $1/n$

## Algorithm

Parameters:  $U, B, \alpha, \delta > 0$

Initialization:  $\theta_{\text{center}} = 0 \in \mathbb{R}^d, t_1 = 1$

For sessions  $i \geq 1$ ,

- Start a new  $EG^\pm$  in  $\mathcal{B}_1(\theta_{\text{center}}, U2^{-i})$  for  $t \geq t_i$
- Get the high probability  $\ell_1$ -ball for  $\theta^*$

$$\|\bar{\theta}_{t-1} - \theta^*\|_1^2 \leq d \|\bar{\theta}_{t-1} - \theta^*\|_2^2 \lesssim \frac{dU2^{-i}B\sqrt{\log(d/\delta)}}{\alpha\sqrt{t-t_i}} =: C(U, t)^2$$

- Define  $t_{i+1} \geq t_i$  as the first time such that  $C(U, t_{i+1}) \leq U2^{-(i+1)}$ , i.e., for

$$\sqrt{t_{i+1} - t_i} \approx \frac{4dB\sqrt{\log(d/\delta)}}{U2^{-i}\alpha}$$

- $\theta_{\text{center}} \leftarrow \bar{\theta}_{t_{i+1}-1}$

After  $n$  time steps, we will have the slow rate high probability bound

$$\mathbb{E}[l_n](\bar{\theta}_{n-1}) - \mathbb{E}[l_n](\hat{\theta}_{t-1}) \lesssim \frac{U2^{-i}B\sqrt{\log(d/\delta)}}{\sqrt{n}} \lesssim \frac{dB^2 \log(d/\delta)}{\alpha n}$$

but with  $U2^{-i} \approx \frac{dB\sqrt{\log(d/\delta)}}{\alpha\sqrt{n}}$ .

# ACCELERATION: FROM SLOW RATE $1/\sqrt{n}$ TO FAST RATE $1/n$

## Algorithm

Parameters:  $U, B, \alpha, \delta > 0$

Initialization:  $\theta_{\text{center}} = 0 \in \mathbb{R}^d, t_1 = 1$

For sessions  $i \geq 1$ ,

- Start a new  $EG^\pm$  in  $\mathcal{B}_1(\theta_{\text{center}}, U2^{-i})$  for  $t \geq t_i$
- Get the high probability  $\ell_1$ -ball for  $\theta^*$

$$\|\bar{\theta}_{t-1} - \theta^*\|_1^2 \leq d \|\bar{\theta}_{t-1} - \theta^*\|_2^2 \lesssim \frac{dU2^{-i}B\sqrt{\log(d/\delta)}}{\alpha\sqrt{t-t_i}} =: C(U, t)^2$$

- Define  $t_{i+1} \geq t_i$  as the first time such that  $C(U, t_{i+1}) \leq U2^{-(i+1)}$ , i.e., for

$$\sqrt{t_{i+1} - t_i} \approx \frac{4dB\sqrt{\log(d/\delta)}}{U2^{-i}\alpha}$$

- $\theta_{\text{center}} \leftarrow \bar{\theta}_{t_{i+1}-1}$

After  $n$  time steps, we will have the slow rate high probability bound

$$\mathbb{E}[l_n](\bar{\theta}_{n-1}) - \mathbb{E}[l_n](\hat{\theta}_{t-1}) \lesssim \frac{U2^{-i}B\sqrt{\log(d/\delta)}}{\sqrt{n}} \lesssim \frac{dB^2 \log(d/\delta)}{\alpha n}$$

but with  $U2^{-i} \approx \frac{dB\sqrt{\log(d/\delta)}}{\alpha\sqrt{n}}$ .

# ACCELERATION: FROM SLOW RATE $1/\sqrt{n}$ TO FAST RATE $1/n$

## Algorithm (SAEW)

Parameters:  $U, B, \alpha, \delta > 0, d_0 \geq 1$

Initialization:  $\theta_{\text{center}} = 0 \in \mathbb{R}^d, t_1 = 1$

For sessions  $i \geq 1$ ,

- Start a new  $EG^\pm$  in  $\mathcal{B}_1(\theta_{\text{center}}, U2^{-i})$  for  $t \geq t_i$
- Get the high probability  $\ell_1$ -ball for  $\theta^*$

$$\|[\bar{\theta}_{t-1}]_{d_0} - \theta^*\|_1^2 \leq d_0 \|[\bar{\theta}_{t-1}]_{d_0} - \theta^*\|_2^2 \lesssim \frac{d_0 U 2^{-i} B \sqrt{\log(d/\delta)}}{\alpha \sqrt{t - t_i}} =: C(U, t)^2$$

where  $[\bar{\theta}_{t-1}]_{d_0}$  is the  $d_0$ -truncation of  $\bar{\theta}_{t-1}$

- Define  $t_{i+1} \geq t_i$  as the first time such that  $C(U, t_{i+1}) \leq U 2^{-(i+1)}$ , i.e., for

$$\sqrt{t_{i+1} - t_i} \approx \frac{4d_0 B \sqrt{\log(d/\delta)}}{U 2^{-i} \alpha}$$

- $\theta_{\text{center}} \leftarrow \bar{\theta}_{t_{i+1}-1}$

After  $n$  time steps, we will have the slow rate high probability bound

$$\mathbb{E}[\ell_n](\bar{\theta}_{n-1}) - \mathbb{E}[\ell_n](\hat{\theta}_{t-1}) \lesssim \frac{U 2^{-i} B \sqrt{\log(d/\delta)}}{\sqrt{n}} \lesssim \frac{d_0 B^2 \log(d/\delta)}{\alpha n}$$

but with  $U 2^{-i} \approx \frac{d_0 B \sqrt{\log(d/\delta)}}{\alpha \sqrt{n}}$ .

# FROM SLOW RATE BOUND TO FAST RATE BOUND : ILLUSTRATION

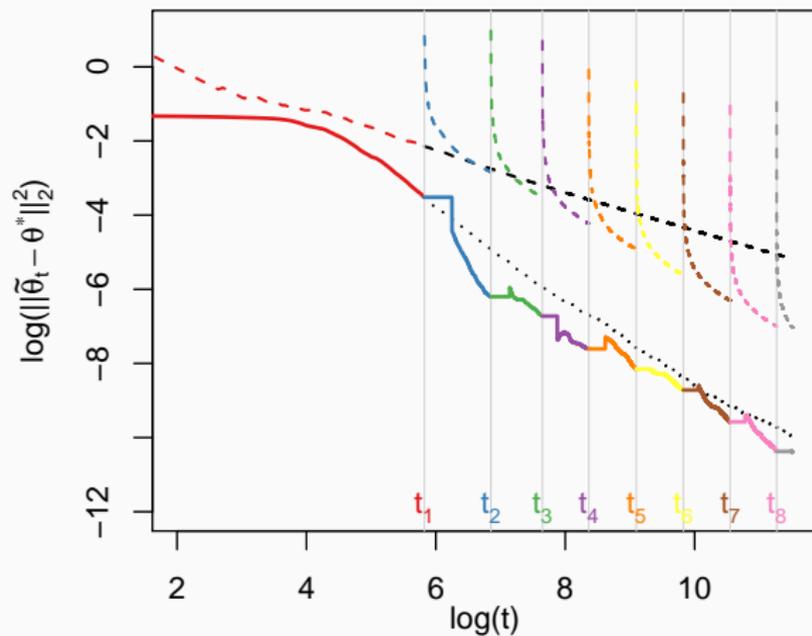


Figure 1: Logarithm of the  $\ell_2$ -error of the averaged estimator.

## Theorem

The average risk of SAEW is upper-bounded as

$$\text{Risk}_{1:n}(\widehat{\theta}_{0:(n-1)}) \lesssim \min \left\{ UB \sqrt{\frac{\log(d/\delta)}{n}}, \frac{d_0 B^2}{\alpha n} \log(d/\delta) \log n + \frac{\alpha U^2}{d_0 n} \right\}.$$

## Remarks:

- Both rates are optimal (in some sense)
- From the strong convexity assumption, this also ensures

$$\|\bar{\theta}_{n-1} - \theta^*\|_2^2 \lesssim \min \left\{ \frac{UB \sqrt{\log(d/\delta)}}{\alpha \sqrt{n}}, \frac{d_0 B^2}{n \alpha^2} \log(d/\delta) \log n + \frac{d_0 U^2}{\alpha n} \right\}$$

- the boundness of the gradients  $B \geq \max_{\theta \in \mathcal{B}_1(0, 2U)} \|\nabla \ell_t(\theta)\|_\infty$  can be weakened to unknown  $B$  under the subgaussian condition

## Theorem

The average risk of SAEW is upper-bounded as

$$\text{Risk}_{1:n}(\hat{\theta}_{0:(n-1)}) \lesssim \min \left\{ UB \sqrt{\frac{\log(d/\delta)}{n}}, \frac{d_0 B^2}{\alpha n} \log(d/\delta) \log n + \frac{\alpha U^2}{d_0 n} \right\}.$$

## Remarks:

- Both rates are optimal (in some sense)
- From the strong convexity assumption, this also ensures

$$\|\tilde{\theta}_{n-1} - \theta^*\|_2^2 \lesssim \min \left\{ \frac{UB \sqrt{\log(d/\delta)}}{\alpha \sqrt{n}}, \frac{d_0 B^2}{n \alpha^2} \log(d/\delta) \log n + \frac{d_0 U^2}{\alpha n^2} \right\}$$

where  $\tilde{\theta}_{n-1} = (t_i - t_{i-1})^{-1} \sum_{t=t_{i-1}}^{t_i-1} \hat{\theta}_{t-1}$  for  $t_i \leq n \leq t_{i+1}$

- the boundness of the gradients  $B \geq \max_{\theta \in \mathcal{B}_1(0, 2U)} \|\nabla \ell_t(\theta)\|_\infty$  can be weakened to unknown  $B$  under the subgaussian condition

We compare three online optimization procedures:

- **RDA**<sup>6</sup>:  $\ell_1$ -regularized dual averaging method

$$\hat{\theta}_t = \arg \min_{\theta \in \mathbb{R}^d} \left\{ \underbrace{\frac{1}{t} \sum_{s=1}^t \nabla \ell_t(\hat{\theta}_{s-1})^\top \theta}_{\text{Linearized loss}} + \underbrace{\lambda \|\theta\|_1}_{\ell_1 \text{ regularization}} + \underbrace{\frac{\gamma}{\sqrt{t}} \|\theta\|_2^2}_{\text{force strong-convexity}} \right\}$$



Good performance for hand-written digits classification

Produces sparse estimators: but slow rate, or fast rate with



No sparse guarantees

- **BOA**<sup>7</sup>: exponential weights with second order regularization ( $\approx EG^\pm$  with good tuning and high probability properties).



achieves fast rate for expert selection



no fast rate in the  $\ell_2$ -ball

- **SAEW**: our acceleration of BOA

All methods are tuned in hindsight with the best parameters on a grid.

<sup>6</sup> Xiao, "Dual averaging methods for regularized stochastic learning and online optimization", 2010

<sup>7</sup> Wintenberger, "Optimal learning with Bernstein Online Aggregation", 2014

Let  $(X_t, Y_t) \in [-X, X]^d \times [-Y, Y]$  be i.i.d. random pairs ( $X, Y > 0$ ).

**Goal:** estimate linearly  $Y_t$  by approaching

$$\theta^* \in \arg \min_{\theta \in \mathbb{R}^d} \mathbb{E}[(Y_t - X_t^\top \theta)^2]$$

The strong convexity assumption is achieved with  $\alpha \leq \lambda_{\min}(\mathbb{E}[X_t X_t^\top])$ .

**Experiment:**  $X_t \sim \mathcal{N}(0, 1)$  for  $d = 500, n = 2000$

$$Y_t = X_t^\top \theta^* + 0.1 \varepsilon_t \quad \text{with } \varepsilon_t \sim \mathcal{N}(0, 1) \text{ i.i.d.}$$

where  $d_0 = \|\theta^*\|_0 = 5, U = \|\theta^*\|_1 = 1$  with non-zero coordinates i.i.d.  $\propto \mathcal{N}(0, 1), \alpha = 1$ .

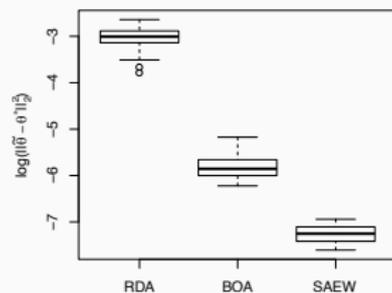


Figure 2: Boxplot (30 simulations)

least square regression with  $d_0 = 5$ ,  $d = 500$ ,  $\sigma = 0.1$

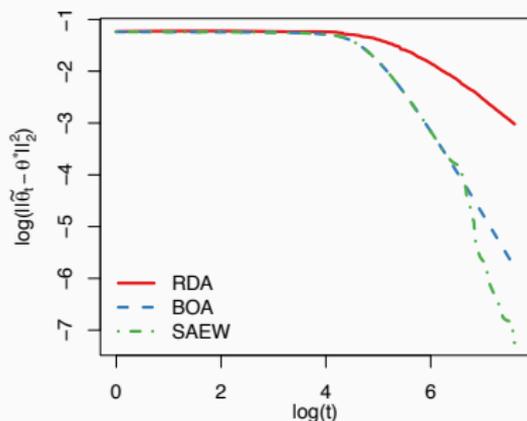


Figure 3: Log of the  $\ell_2$  error.

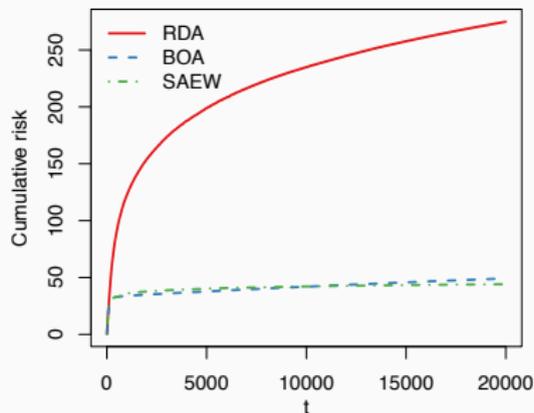


Figure 4: Cumulative risk.

**Remarks:** the cumulative risks are at most of order:

$$\text{RDA: } \sigma^2 d \log n$$

$$\text{BOA: } \sigma^2 \sqrt{n \log d} + \log d$$

$$\text{SAEW: } \sigma^2 d_0 \log d \log n + \log d$$

**In practice:** much better performance if we allow multiple pass on the training set.

Simulation: least square regression with  $d_0 = 2$ ,  $d = 2$ ,  $\sigma = 0.3$

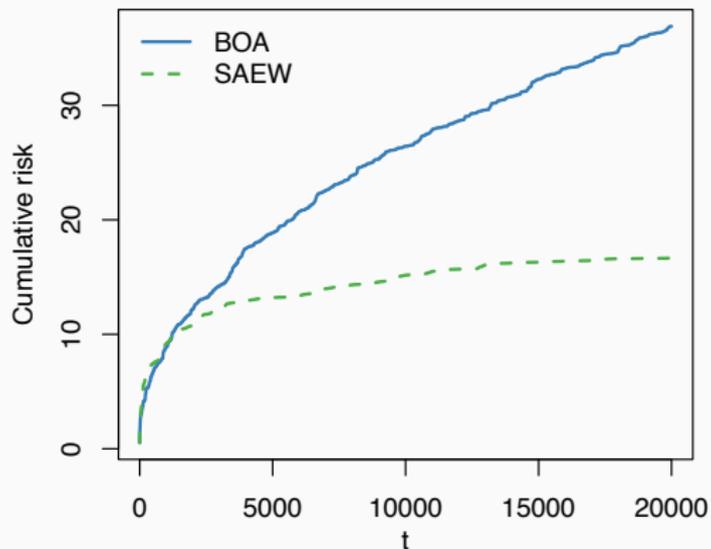


Figure 5: Cumulative risk for square linear regression with  $d = d_0 = 2$ .

Let  $X, Y > 0$ . Let  $(X_t, Y_t) \in [-X, X]^d \times [-Y, Y]$  be i.i.d. random pairs.

### Theorem

SAEW applied with  $B = 2X(Y + 2XU)$  satisfies

$$\text{Risk}_{1:n}(\widehat{\theta}_{0:(n-1)}) \lesssim \min \left\{ UX(Y + XU) \sqrt{\frac{\log(d/\delta)}{n}}, \frac{d_0 X^2 (Y^2 + X^2 Y^2)}{\alpha n} \log \frac{d}{\delta} \log n + \frac{\alpha U^2}{d_0 n} \right\}.$$

A better tuning of  $EG^\pm$  with  $\eta_t \approx 1/\sqrt{\sum_{s=1}^t \|\nabla \ell_s(\widehat{\theta}_{s-1})\|_\infty^2}$  allows to substituted  $B^2$  with  $X^2 \sigma^2$  with

$$\sigma^2 = \mathbb{E}[(Y_t - X_t^\top \theta^*)^2].$$

in the instantaneous risk of  $\tilde{\theta}_{n-1}$ .

Let  $X, Y > 0$ . Let  $(X_t, Y_t) \in [-X, X]^d \times [-Y, Y]$  be i.i.d. random pairs.

### Theorem

SAEW applied with  $B = 2X(Y + 2XU)$  satisfies

$$\text{Risk}(\tilde{\theta}_{n-1}) \lesssim \min \left\{ \sigma \sqrt{\frac{\log(d/\delta)}{n}}, \frac{d_0 X^2 \sigma^2}{n\alpha} \log \frac{d}{\delta} + \frac{\alpha U^2}{d_0 n^2} \right\}.$$

A better tuning of  $EG^\pm$  with  $\eta_t \approx 1/\sqrt{\sum_{s=1}^t \|\nabla \ell_s(\hat{\theta}_{s-1})\|_\infty^2}$  allows to substituted  $B^2$  with  $X^2 \sigma^2$  with

$$\sigma^2 = \mathbb{E}[(Y_t - X_t^\top \theta^*)^2].$$

in the instantaneous risk of  $\tilde{\theta}_{n-1}$ .

👍 Optimal rate for sparse least square regression.

The algorithm needs to know :  $d_0, U, B, \alpha, \delta$

💡 Run a meta-algorithm (BOA<sup>8</sup>) with parameters in a growing grid.

- 👎 We leave the initial setting since we need
  - to observe the gradients of all sub-algorithms.
  - to clip the predictions  $\widehat{\theta}_{t-1}^\top X_t \rightarrow [\widehat{\theta}_{t-1}^\top X_t]_{[-Y, Y]}$  otherwise we pay the maximal value of  $U$  considered in the final bound.
  - we need strongly convex  $\ell_t$  (instead of  $\mathbb{E}[\ell_t]$  only)
- 👍 This works to build an estimator for least square regression.

## Theorem (Calibrated SAEW for Least Square Linear Regression)

The excess risk of the estimator produced by the meta-algorithm is of order

$$\mathcal{O}_n \left( \underbrace{\frac{Y^2}{n} \log \left( \frac{(\log d)(\log n + \log Y)}{\delta} \right)}_{\text{Price of calibration}} + \underbrace{\frac{d_0 X^2 \sigma^2}{\alpha^* n} \log(d/\delta)}_{\text{Fast rate with best } \alpha^*, d_0} \right),$$

where  $\alpha^*$  is the largest strong-convexity parameter.

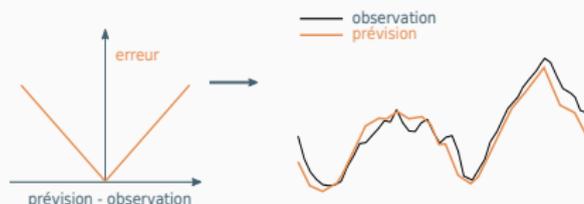
Can we substitute  $\alpha^*$  with local strong convexity (cf. Lasso)?



<sup>8</sup> Wintenberger, "Optimal learning with Bernstein Online Aggregation", 2014.

Let  $\tau \in (0, 1)$ . Let  $(X_t, Y_t) \in \mathbb{R}^d \times \mathbb{R}$  be i.i.d. random pairs.

**Goal:** estimate the conditional  $\tau$ -quantile of  $Y_t$  given  $X_t$ .



**Popular solution:** linear regression with the pinball loss by  $\rho_\tau : u \in \mathbb{R} \rightarrow u(\tau - \mathbf{1}_{u < 0})$   
 The conditional quantile  $q_\tau(Y_t|X_t)$  is the solution of

$$q_\tau(Y_t|X_t) \in \arg \min_g \mathbb{E}[\rho_\tau(Y_t - g(X_t)) | X_t].$$

→ minimize the pinball loss  $q_\tau = \tau$ -quantile prediction



non-strongly convex loss

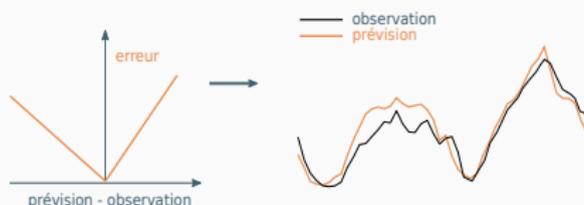


strongly convex risk<sup>9</sup> → ok for fixed parameter

<sup>9</sup> Steinwart and Christmann, "Estimating conditional quantiles with the help of the pinball loss", 2011.

Let  $\tau \in (0, 1)$ . Let  $(X_t, Y_t) \in \mathbb{R}^d \times \mathbb{R}$  be i.i.d. random pairs.

**Goal:** estimate the conditional  $\tau$ -quantile of  $Y_t$  given  $X_t$ .



**Popular solution:** linear regression with the pinball loss by  $\rho_\tau : u \in \mathbb{R} \rightarrow u(\tau - \mathbf{1}_{u < 0})$   
 The conditional quantile  $q_\tau(Y_t|X_t)$  is the solution of

$$q_\tau(Y_t|X_t) \in \arg \min_g \mathbb{E}[\rho_\tau(Y_t - g(X_t)) | X_t].$$

→ minimize the pinball loss  $q_\tau = \tau$ -quantile prediction



non-strongly convex loss

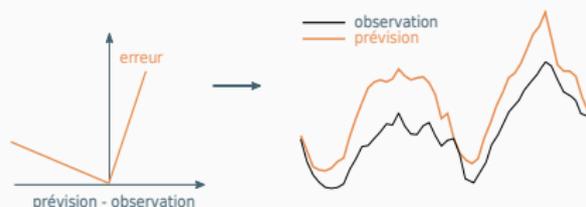


strongly convex risk<sup>9</sup> → ok for fixed parameter

<sup>9</sup> Steinwart and Christmann, "Estimating conditional quantiles with the help of the pinball loss", 2011.

Let  $\tau \in (0, 1)$ . Let  $(X_t, Y_t) \in \mathbb{R}^d \times \mathbb{R}$  be i.i.d. random pairs.

**Goal:** estimate the conditional  $\tau$ -quantile of  $Y_t$  given  $X_t$ .



**Popular solution:** linear regression with the pinball loss by  $\rho_\tau : u \in \mathbb{R} \rightarrow u(\tau - \mathbb{1}_{u < 0})$   
 The conditional quantile  $q_\tau(Y_t|X_t)$  is the solution of

$$q_\tau(Y_t|X_t) \in \arg \min_g \mathbb{E}[\rho_\tau(Y_t - g(X_t)) | X_t].$$

→ minimize the pinball loss  $q_\tau = \tau$ -quantile prediction



non-strongly convex loss

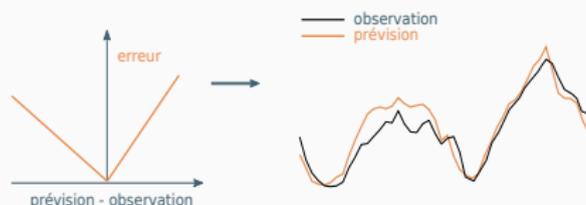


strongly convex risk<sup>9</sup> → ok for fixed parameter

<sup>9</sup> Steinwart and Christmann, "Estimating conditional quantiles with the help of the pinball loss", 2011.

Let  $\tau \in (0, 1)$ . Let  $(X_t, Y_t) \in \mathbb{R}^d \times \mathbb{R}$  be i.i.d. random pairs.

**Goal:** estimate the conditional  $\tau$ -quantile of  $Y_t$  given  $X_t$ .



**Popular solution:** linear regression with the pinball loss by  $\rho_\tau : u \in \mathbb{R} \rightarrow u(\tau - \mathbf{1}_{u < 0})$   
 The conditional quantile  $q_\tau(Y_t|X_t)$  is the solution of

$$q_\tau(Y_t|X_t) \in \arg \min_g \mathbb{E}[\rho_\tau(Y_t - g(X_t)) | X_t].$$

→ minimize the pinball loss  $q_\tau = \tau$ -quantile prediction



non-strongly convex loss

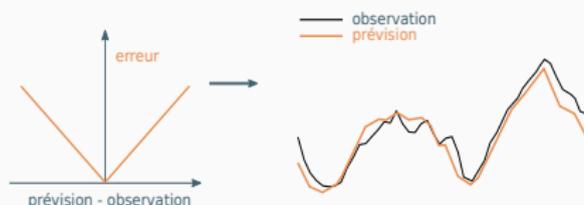


strongly convex risk<sup>9</sup> → ok for fixed parameter

<sup>9</sup> Steinwart and Christmann, "Estimating conditional quantiles with the help of the pinball loss", 2011.

Let  $\tau \in (0, 1)$ . Let  $(X_t, Y_t) \in \mathbb{R}^d \times \mathbb{R}$  be i.i.d. random pairs.

**Goal:** estimate the conditional  $\tau$ -quantile of  $Y_t$  given  $X_t$ .



**Popular solution:** linear regression with the pinball loss by  $\rho_\tau : u \in \mathbb{R} \rightarrow u(\tau - \mathbf{1}_{u < 0})$   
 The conditional quantile  $q_\tau(Y_t|X_t)$  is the solution of

$$q_\tau(Y_t|X_t) \in \arg \min_g \mathbb{E}[\rho_\tau(Y_t - g(X_t)) | X_t].$$

→ minimize the pinball loss  $q_\tau = \tau$ -quantile prediction



non-strongly convex loss

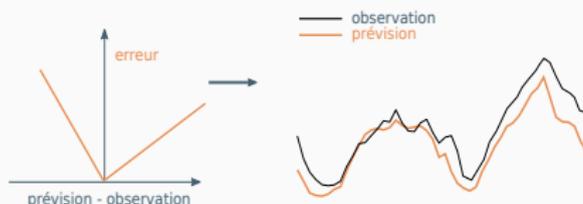


strongly convex risk<sup>9</sup> → ok for fixed parameter

<sup>9</sup> Steinwart and Christmann, "Estimating conditional quantiles with the help of the pinball loss", 2011.

Let  $\tau \in (0, 1)$ . Let  $(X_t, Y_t) \in \mathbb{R}^d \times \mathbb{R}$  be i.i.d. random pairs.

**Goal:** estimate the conditional  $\tau$ -quantile of  $Y_t$  given  $X_t$ .



**Popular solution:** linear regression with the pinball loss by  $\rho_\tau : u \in \mathbb{R} \rightarrow u(\tau - \mathbf{1}_{u < 0})$   
 The conditional quantile  $q_\tau(Y_t|X_t)$  is the solution of

$$q_\tau(Y_t|X_t) \in \arg \min_g \mathbb{E}[\rho_\tau(Y_t - g(X_t)) | X_t].$$

→ minimize the pinball loss  $q_\tau = \tau$ -quantile prediction



non-strongly convex loss

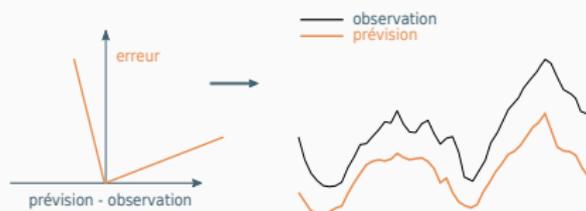


strongly convex risk<sup>9</sup> → ok for fixed parameter

<sup>9</sup> Steinwart and Christmann, "Estimating conditional quantiles with the help of the pinball loss", 2011.

Let  $\tau \in (0, 1)$ . Let  $(X_t, Y_t) \in \mathbb{R}^d \times \mathbb{R}$  be i.i.d. random pairs.

**Goal:** estimate the conditional  $\tau$ -quantile of  $Y_t$  given  $X_t$ .



**Popular solution:** linear regression with the pinball loss by  $\rho_\tau : u \in \mathbb{R} \rightarrow u(\tau - \mathbf{1}_{u < 0})$   
 The conditional quantile  $q_\tau(Y_t|X_t)$  is the solution of

$$q_\tau(Y_t|X_t) \in \arg \min_g \mathbb{E}[\rho_\tau(Y_t - g(X_t)) | X_t].$$

→ minimize the pinball loss  $q_\tau = \tau$ -quantile prediction



non-strongly convex loss

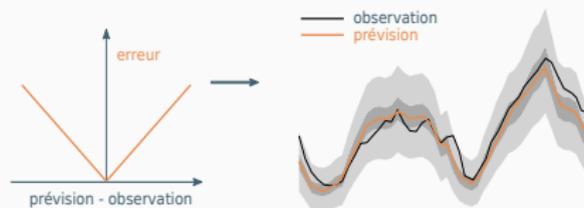


strongly convex risk<sup>9</sup> → ok for fixed parameter

<sup>9</sup> Steinwart and Christmann, "Estimating conditional quantiles with the help of the pinball loss", 2011.

Let  $\tau \in (0, 1)$ . Let  $(X_t, Y_t) \in \mathbb{R}^d \times \mathbb{R}$  be i.i.d. random pairs.

**Goal:** estimate the conditional  $\tau$ -quantile of  $Y_t$  given  $X_t$ .



**Popular solution:** linear regression with the pinball loss by  $\rho_\tau : u \in \mathbb{R} \rightarrow u(\tau - \mathbb{1}_{u < 0})$   
 The conditional quantile  $q_\tau(Y_t|X_t)$  is the solution of

$$q_\tau(Y_t|X_t) \in \arg \min_g \mathbb{E}[\rho_\tau(Y_t - g(X_t)) | X_t].$$

→ minimize the pinball loss  $q_\tau = \tau$ -quantile prediction



non-strongly convex loss



strongly convex risk<sup>9</sup> → ok for fixed parameter

<sup>9</sup> Steinwart and Christmann, "Estimating conditional quantiles with the help of the pinball loss", 2011.

Experiment:  $d_0 = 5, d = 100$

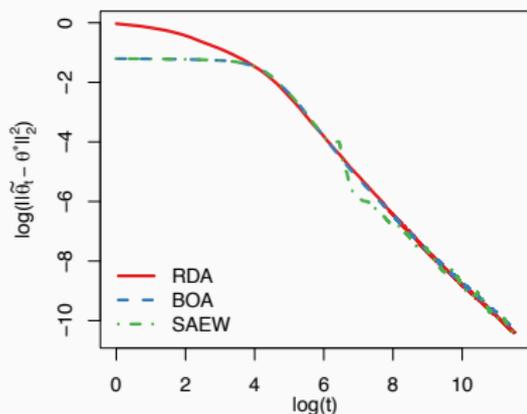


Figure 6: Log. of the  $\ell_2$ -error of  $\bar{\theta}_t$

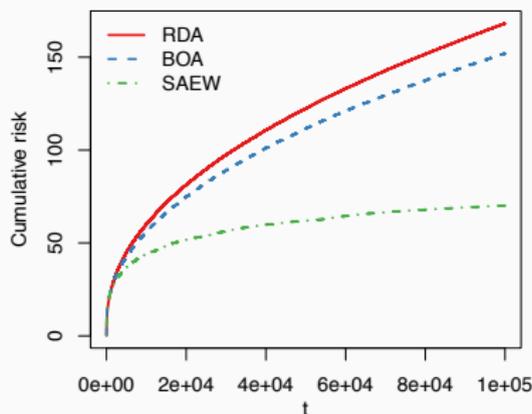


Figure 7: Cumulative risk

All the methods empirically get the fast rate  $1/n$  for the  $\ell_2$ -error of the estimator...

**But** only SAEW

- has the theoretical guarantee

- has  $\mathcal{O}(\log n)$  cumulative risk



Some future work:

- Is averaging an efficient acceleration procedure for  $EG^\pm$ ?
- Calibration of the parameters in the original online optimization setting
- Produce sparse estimators  $\hat{\theta}_{t-1}$ 
  - improve the dependency on the strong convexity parameter (only local)
- Oracle bound: no assumption on the sparsity of  $\theta^*$

THANK YOU !

## REFERENCES

---

-  Bunea, F., A. Tsybakov, and M. Wegkamp. “Aggregation for Gaussian regression”. In: *The Annals of Statistics* 35.4 (2007), pp. 1674–1697.
-  Cesa-Bianchi, N. and G. Lugosi. *Prediction, Learning, and Games*. Cambridge University Press, 2006.
-  Gaillard, P. and O. Wintenberger. “Sparse Accelerated Exponential Weights”. Accepted at AISTAT’17. 2017.
-  Gerchinovitz, S. “Sparsity regret bounds for individual sequences in online linear regression”. In: *The Journal of Machine Learning Research* 14.1 (2013), pp. 729–769.
-  Kivinen, J. and M. K. Warmuth. “Exponentiated Gradient Versus Gradient Descent for Linear Predictors”. In: *Information and Computation* 132.1 (1997), pp. 1–63.
-  Rigollet, P. and A. Tsybakov. “Exponential screening and optimal rates of sparse estimation”. In: *The Annals of Statistics* (2011), pp. 731–771.
-  Steinwart, I. and A. Christmann. “Estimating conditional quantiles with the help of the pinball loss”. In: *Bernoulli* 17.1 (2011), pp. 211–225.
-  Wintenberger, O. “Optimal learning with Bernstein Online Aggregation”. In: Extended version available at arXiv:1404.1356 [stat. ML] (2014).
-  Xiao, L. “Dual averaging methods for regularized stochastic learning and online optimization”. In: *Journal of Machine Learning Research* 11 (2010), pp. 2543–2596.