

ACCÉLÉRATION PARCIMONIEUSE DES POIDS EXPONENTIELS

Pierre Gaillard

30 mai 2017

INRIA Paris – SIERRA

Travail en collaboration avec Olivier Wintenberger

CADRE

Minimisation **séquentielle** d'une suite i.i.d. de fonctions de pertes¹ $\ell_1, \dots, \ell_n : \mathbb{R}^d \rightarrow \mathbb{R}$.

Cadre : À chaque pas de temps $t = 1, \dots, n$

- le statisticien met à jour son estimateur $\hat{\theta}_{t-1} \in \mathbb{R}^d$ basé sur le passé $\ell_1, \dots, \ell_{t-1}$
- on observe ℓ_t .

Objectif : minimiser le risque de $\hat{\theta}_n$

$$\text{Risk}(\hat{\theta}_n) := \mathbb{E}[\ell_{n+1}](\hat{\theta}_n).$$

Exemple : observation séquentielle d'un échantillon $\{(X_t, Y_t)\}_{1 \leq t \leq n}$ avec la perte carrée : $\ell_t : \theta \mapsto (Y_t - X_t \cdot \theta)^2$.

$$\text{Risk}(\theta) = \mathbb{E}[(Y_t - X_t \cdot \theta)^2].$$

1. CESA-BIANCHI et LUGOSI, *Prediction, Learning, and Games*, 2006.

On suppose que le domaine, les gradients et pertes sont bornées ps. (ou sous-gaussien).

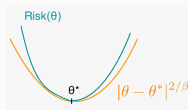
Hypothèse de Bernstein (ou Lojasiewicz's)

$$\exists \mu > 0, 0 \leq \beta \leq 1, \theta^* \in \mathbb{R}^d, \forall \theta \in \mathbb{R}^d \quad \mu \|\theta - \theta^*\|_2^2 \leq [\text{Risk}(\theta) - \text{Risk}(\theta^*)]^\beta.$$

Remarques :

- **Risque fortement convexe** \rightarrow vérifié avec $\beta = 1$
 \rightarrow peut être vérifié pour des pertes non fortement convexes comme la perte absolue ou quantile, si les distributions sont assez régulières.
- **Risque seulement convexe** \rightarrow vérifié avec $\mu = 1/(\text{Borne sur le domaine})$ et $\beta = 0$
- Tout un panel entre les deux même dans des cadres qui semblent discrets comme la classification.

Ex : $Y \in \{-1, 1\}$, perte Hinge^a \rightarrow vérifié avec $\beta > 0$



a. LECUÉ, "Optimal Rates of Aggregation in Classification Under Low Noise Assumption", 2006.

On suppose que le domaine, les gradients et pertes sont bornées ps. (ou sous-gaussien).

Hypothèse de Bernstein (ou Lojasiewicz's)

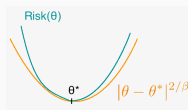
$$\exists \mu > 0, 0 \leq \beta \leq 1, \theta^* \in \mathbb{R}^d, \forall \theta \in \mathbb{R}^d \quad \mu \|\theta - \theta^*\|_2^2 \leq [\text{Risk}(\theta) - \text{Risk}(\theta^*)]^\beta .$$

$$\theta^* \text{ is } d_0\text{-sparse} \quad \|\theta^*\|_1 \leq U$$

Remarques :

- **Risque fortement convexe** \rightarrow vérifié avec $\beta = 1$
 \rightarrow peut être vérifié pour des pertes non fortement convexes comme la perte absolue ou quantile, si les distributions sont assez régulières.
- **Risque seulement convexe** \rightarrow vérifié avec $\mu = 1/(\text{Borne sur le domaine})$ et $\beta = 0$
- Tout un panel entre les deux même dans des cadres qui semblent discrets comme la classification.

Ex : $Y \in \{-1, 1\}$, perte Hinge^a \rightarrow vérifié avec $\beta > 0$



a. LECUÉ, "Optimal Rates of Aggregation in Classification Under Low Noise Assumption", 2006.

On suppose que le domaine, les gradients et pertes sont bornées ps. (ou sous-gaussien).

Hypothèse de Bernstein (ou Lojasiewicz's)

$$\exists \mu > 0, 0 \leq \beta \leq 1, \theta^* \in \mathbb{R}^d, \forall \theta \in \mathbb{R}^d \quad \mu \|\theta - \theta^*\|_2^2 \leq [\text{Risk}(\theta) - \text{Risk}(\theta^*)]^\beta .$$

θ^* is d_0 -sparse $\|\theta^*\|_1 \leq U$

Résultat : on propose une procédure qui redémarre des algorithmes de minimisation dans la boule ℓ_1 pour obtenir de meilleures vitesses (en n)

$$\text{Risk}(\hat{\theta}_n) \lesssim \text{Risk}(\theta^*) + \min \left\{ \left(\frac{d_0 \log d}{\mu n} \right)^{\frac{1}{2-\beta}}, U \sqrt{\frac{\log d}{n}} \right\} .$$

Vitesse rapide : mieux pour n, μ grands

Vitesse lente : mieux pour n, μ petits

Procédure	Sequentielle	Vitesse ($\beta = 1$)	Polynomiale	Hypothèse
Lasso ²	✗	$\frac{d_0 \log d}{\mu n}$	✓	Bernstein ³
EWA + sparsity patern ⁴	✗	$\frac{d_0 \log d}{n}$	✗	Forte Conv.
SeqSEW ⁵	✓	$\frac{d_0 \log d}{n}$	✗	Forte Conv.
ℓ_1 -RDA method ⁶	✓	$\frac{d}{n}$	✓	Forte Conv.
Agarwal ⁷ , Steinhard ⁸	✓	$\frac{d_0 \log d}{\mu n}$	✓	Forte Conv.
SAEW	✓	$\frac{d_0 \log d}{\mu n}$	✓	Bernstein

2. BUNEA, TSYBAKOV et WEGKAMP, "Aggregation for Gaussian regression", 2007.

3. PIERRE, VINCENT et GUILLAUME, "Estimation bounds and sharp oracle inequalities of regularized procedures with Lipschitz loss functions", 2017.

4. RIGOLLET et TSYBAKOV, "Exponential screening and optimal rates of sparse estimation", 2011.

5. GERCHINOVITZ, "Sparsity regret bounds for individual sequences in online linear regression", 2013.

6. XIAO, "Dual averaging methods for regularized stochastic learning and online optimization", 2010.

7. AGARWAL, NEGAHBAN et WAINWRIGHT, "Stochastic optimization and sparse statistical recovery : Optimal algorithms for high dimensions", 2012.

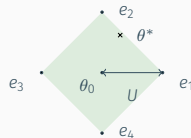
8. STEINHARDT, WAGER et LIANG, "The Statistics of Streaming Sparse Regression", 2014.

OPTIMISATION CONVEXE DANS LA BOULE
 l_1 -BALL AVEC UNE VITESSE LENTE

Obectif : approcher θ^* dans la boule ℓ_1 :

$$\mathcal{B}_1(\theta_0, U) := \{\theta \in \mathbb{R}^d : \|\theta - \theta_0\|_1 \leq U\}$$

On définit $e_k := \theta_0 + (0, \dots, 0, \pm U, 0, \dots, 0)$ pour $k = 1, \dots, 2d$



L'algorithme EG^\pm

Kivinen et Warmuth, 1997

Paramètres : tuning parameters $\eta_t > 0$.

pour chaque itération $t = 1, 2, \dots$ **faire**

- Définir $\theta_{t-1} = \sum_{k=1}^{2d} p_{k,t} e_k$ et prévoir $\hat{\theta}_{t-1} = \sum_{s=1}^{t-1} \theta_s$.
- Mettre à jour les poids pour $k = 1, \dots, 2d$

$$p_{k,t} = \frac{e^{-\eta_t \sum_{s=1}^t \nabla \ell_s(\theta_{s-1}) \cdot e_k}}{\sum_{j=1}^{2d} e^{-\eta_t \sum_{s=1}^t \nabla \ell_s(\theta_{s-1}) \cdot e_j}}$$

fin

Performance : le choix $\eta_t \approx 1/(U\sqrt{t})$ assure avec grande probabilité

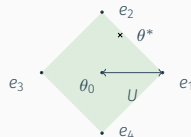
$$\text{Risk}(\hat{\theta}_n) - \text{Risk}(\theta^*) \lesssim U \sqrt{\frac{\log d}{n}}.$$

Remarque : beaucoup d'autres sous-routines peuvent être utilisées !

Obectif : approcher θ^* dans la boule ℓ_1 :

$$\mathcal{B}_1(\theta_0, U) := \{\theta \in \mathbb{R}^d : \|\theta - \theta_0\|_1 \leq U\}$$

On définit $e_k := \theta_0 + (0, \dots, 0, \pm U, 0, \dots, 0)$ pour $k = 1, \dots, 2d$



L'algorithme EG^\pm

Kivinen et Warmuth, 1997

Paramètres : tuning parameters $\eta_t > 0$.

pour chaque itération $t = 1, 2, \dots$ faire

- Définir $\theta_{t-1} = \sum_{k=1}^{2d} p_{k,t} e_k$ et prévoir $\hat{\theta}_{t-1} = \sum_{s=1}^{t-1} \theta_s$.
- Mettre à jour les poids pour $k = 1, \dots, 2d$

$$p_{k,t} = \frac{e^{-\eta_t \sum_{s=1}^t \nabla \ell_s(\theta_{s-1}) \cdot e_k}}{\sum_{j=1}^{2d} e^{-\eta_t \sum_{s=1}^t \nabla \ell_s(\theta_{s-1}) \cdot e_j}}$$

fin

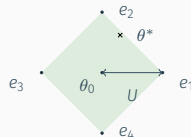
Performance : le choix $\eta_t \approx 1/(U\sqrt{t})$ assure avec grande probabilité

$$\mu \|\hat{\theta}_n - \theta^*\|_2^{2/\beta} \leq \text{Risk}(\hat{\theta}_n) - \text{Risk}(\theta^*) \lesssim U \sqrt{\frac{\log d}{n}}.$$

Obectif : approcher θ^* dans la boule ℓ_1 :

$$\mathcal{B}_1(\theta_0, U) := \{\theta \in \mathbb{R}^d : \|\theta - \theta_0\|_1 \leq U\}$$

On définit $e_k := \theta_0 + (0, \dots, 0, \pm U, 0, \dots, 0)$ pour $k = 1, \dots, 2d$



L'algorithme EG^\pm

Kivinen et Warmuth, 1997

Paramètres : tuning parameters $\eta_t > 0$.

pour chaque itération $t = 1, 2, \dots$ **faire**

- Définir $\theta_{t-1} = \sum_{k=1}^{2d} p_{k,t} e_k$ et prévoir $\hat{\theta}_{t-1} = \sum_{s=1}^{t-1} \theta_s$.
- Mettre à jour les poids pour $k = 1, \dots, 2d$

$$p_{k,t} = \frac{e^{-\eta_t \sum_{s=1}^t \nabla \ell_s(\theta_{s-1}) \cdot e_k}}{\sum_{j=1}^{2d} e^{-\eta_t \sum_{s=1}^t \nabla \ell_s(\theta_{s-1}) \cdot e_j}}$$

fin

Performance : le choix $\eta_t \approx 1/(U\sqrt{t})$ assure avec grande probabilité

$$\|\hat{\theta}_n - \theta^*\|_2 \lesssim \frac{1}{\sqrt{\mu}} \left(U \sqrt{\frac{\log d}{n}} \right)^{\frac{\beta}{2}} . \quad (*)$$

Lemma (Hoeffding)

Si $-B \leq X \leq B$ est une variable aléatoire. Alors,

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

1. On borne les gradients instantanés

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

2. On somme sur t , la somme télescope

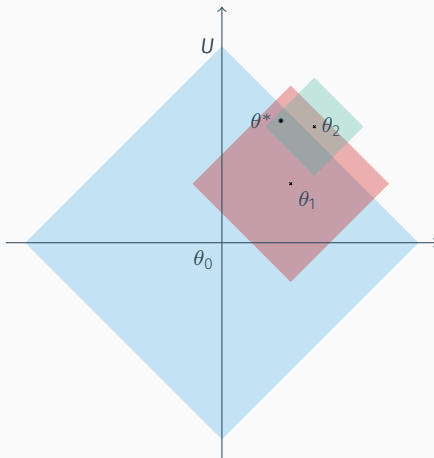
$$\begin{aligned} \sum_{t=1}^n \ell_t(\theta_{t-1}) - \ell_t(\theta^*) &\stackrel{\text{Jensen}}{\leq} \sum_1^n \nabla \ell_t(\theta_{t-1})^\top (\theta_{t-1} - \theta^*) \leq \max_{k=1, \dots, 2d} \left\{ \sum_1^n \nabla \ell_t(\theta_{t-1})^\top (\theta_{t-1} - e_k) \right\} \\ &\leq \max_{k=1, \dots, 2d} \left\{ \frac{1}{\eta} \log \frac{p_{k,n}}{p_{k,0}} + \frac{\eta B n}{4} \right\} \leq \frac{\log(2d)}{\eta} + \frac{\eta B n}{4} \end{aligned}$$

3. On majore le risque par le regret avec les inégalités d'Hoeffding et de Jensen

ACCÉLÉRATION : DE LA VITESSE LENTE À
LA VITESSE RAPIDE

ACCELERATION : DE LA VITESSE LENTE À LA VITESSE RAPIDE

Idée : redémarrer régulièrement la sous-routine dans des boules de plus en plus petites autour de l'estimateur courant.



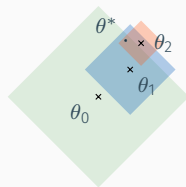
ACCELERATION : DE LA VITESSE LENTE À LA VITESSE RAPIDE

Pour $i = 1, \dots, i_n$, **redémarrer** la sous-routine pour optimiser dans des boules de plus en plus petites de rayon $U_i := U2^{-i}$.

Où ? Centrées sur l'estimateur courant $\hat{\theta}_t$.

Quand ? Quand on sait qu'avec grande probabilité

$$\|\hat{\theta}_t - \theta^*\|_1 \leq U_i.$$



Calcul du quand (attention il faut un peu jongler) : d'après la vitesse lente :

$$\|\hat{\theta}_{t_i} - \theta^*\|_1 \leq \sqrt{d} \|\hat{\theta}_{t_i} - \theta^*\|_2 \stackrel{(*)}{\lesssim} \sqrt{\frac{d}{\mu}} \left(U_{i-1} \sqrt{\frac{\log d}{t_i}} \right)^{\frac{\beta}{2}}.$$

Comme $U_{i-1} = 2U_i$, il suffit donc que

$$\sqrt{\frac{d}{\mu}} \left(2U_i \sqrt{\frac{\log d}{t_i}} \right)^{\frac{\beta}{2}} \lesssim U_i \Leftrightarrow t_i \approx \left(\frac{d}{\mu U_i^{2-\beta}} \right)^{\frac{2}{\beta}} \log d$$

L'algorithme : redémarrer la sous-routine en $\hat{\theta}_t$ tous les

$$t_i \approx \left(\frac{d}{\mu U_i^{2-\beta}} \right)^{\frac{2}{\beta}} \log d$$

instants.

Que se passe-t-il ? La sous-routine travaille dans des boules de plus en plus petites,

$$U_i \approx \left(\frac{d \log d}{\mu t_i} \right)^{\frac{1}{2-\beta}} \sqrt{\frac{t_i}{\log d}}$$

L'algorithme : redémarrer la sous-routine en $\hat{\theta}_t$ tous les

$$t_i \approx \left(\frac{d}{\mu U_i^{2-\beta}} \right)^{\frac{2}{\beta}} \log d$$

instants.

Que se passe t-il ? La sous-routine travaille dans des boules de plus en plus petites,

$$U_{i_n} \approx \left(\frac{d \log d}{\mu n} \right)^{\frac{1}{2-\beta}} \sqrt{\frac{n}{\log d}}$$

Transformation de la vitesse lente en vitesse rapide : la vitesse lente de la dernière sous-routine (au temps n) devient

$$\text{Risk}(\hat{\theta}_n) - \text{Risk}(\theta^*) \lesssim U_{i_n} \sqrt{\frac{\log d}{n}} \lesssim \left(\frac{d \log d}{\mu n} \right)^{\frac{1}{2-\beta}}$$

L'algorithme : redémarrer la sous-routine en $\hat{\theta}_t$ tous les

$$t_i \approx \left(\frac{d}{\mu U_i^{2-\beta}} \right)^{\frac{2}{\beta}} \log d$$

instants.

Que se passe t-il ? La sous-routine travaille dans des boules de plus en plus petites,

$$U_{i_n} \approx \left(\frac{d \log d}{\mu n} \right)^{\frac{1}{2-\beta}} \sqrt{\frac{n}{\log d}}$$

Transformation de la vitesse lente en vitesse rapide : la vitesse lente de la dernière sous-routine (au temps n) devient

$$\text{Risk}(\hat{\theta}_n) - \text{Risk}(\theta^*) \lesssim U_{i_n} \sqrt{\frac{\log d}{n}} \lesssim \left(\frac{d \log d}{\mu n} \right)^{\frac{1}{2-\beta}}$$

Mais on ne voulait pas d_0 plutôt que d ?

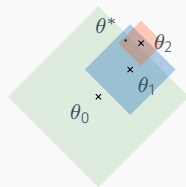
ACCELERATION : DE LA VITESSE LENTE À LA VITESSE RAPIDE

Pour $i = 1, \dots, i_n$, **redémarrer** la sous-routine pour optimiser dans des boules de plus en plus petites de rayon $U_i := U2^{-i}$.

Où ? Centrées sur l'estimateur courant tronqué $[\hat{\theta}_t]_{d_0}$ à d_0 coordonnées non-nulles.

Quand ? Quand on sait qu'avec grande probabilité

$$\|[\hat{\theta}_t]_{d_0} - \theta^*\|_1 \leq U_i.$$



Calcul du quand (attention il faut un peu jongler) : d'après la vitesse lente :

$$\|\hat{\theta}_{t_i} - \theta^*\|_1 \leq \sqrt{d_0} \|\hat{\theta}_{t_i} - \theta^*\|_2 \stackrel{(*)}{\lesssim} \sqrt{\frac{d_0}{\mu}} \left(U_{i-1} \sqrt{\frac{\log d}{t_i}} \right)^{\frac{\beta}{2}}.$$

Les mêmes calculs [...] donnent

$$\text{Risk}(\hat{\theta}_n) - \text{Risk}(\theta^*) \lesssim U_{i_n} \sqrt{\frac{\log d}{n}} \lesssim \left(\frac{d_0 \log d}{\mu n} \right)^{\frac{1}{2-\beta}}$$

VITESSE LENTE À VITESSE RAPIDE : ILLUSTRATION

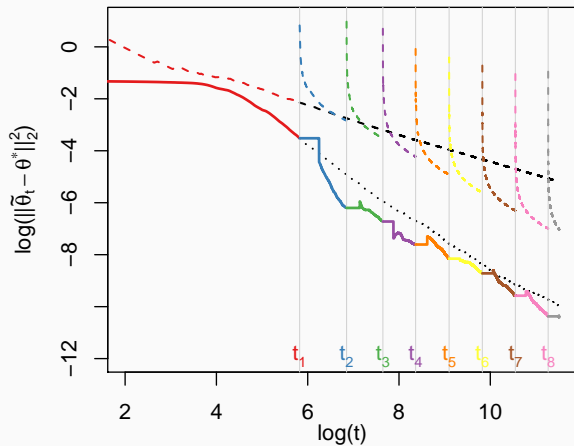


FIGURE 1 – Logarithm of the ℓ_2 -error of the averaged estimator.

Échantillon $(X_t, Y_t) \in [-X, X]^d \times [-Y, Y]$ i.i.d.

Objectif : estimer linéairement Y_t en approchant

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

Le risque est fortement convexe $\beta = 1$ avec $\mu \leq \lambda_{\min}(\mathbb{E}[X_t X_t^\top])$.

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

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

où $d_0 = \|\theta^*\|_0 = 5, U = \|\theta^*\|_1 = 1$, coordonnées non nulles i.i.d. $\propto \mathcal{N}(0, 1), \mu = 1$.

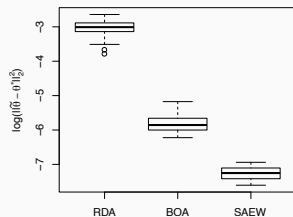


FIGURE 2 – Boxplot (30 simulations)

Simulations : $d_0 = 5$, $d = 500$, $\sigma = 0.1$

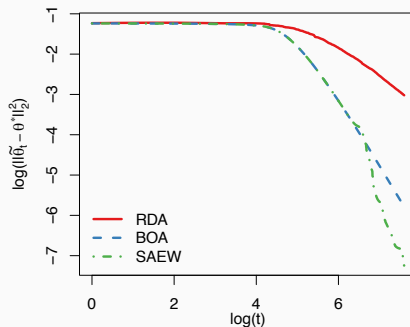


FIGURE 3 – Log de l'erreur ℓ_2 .

En pratique : bien meilleur si on autorise plusieurs passes sur les données.



Travail en court et futur :

- Calibration des paramètres plus propre (sans surcouche) : procédure complètement adaptative.
- Borne oracle : pas d'hypothèse de parcimonie sur θ^*
- Le faire sans recommencer ?
- Applications...

MERCI !

RÉFÉRENCES



AGARWAL, A., S. NEGAHBAN et M. J. WAINWRIGHT. "Stochastic optimization and sparse statistical recovery : Optimal algorithms for high dimensions". In : *Advances in Neural Information Processing Systems*. 2012, p. 1538–1546.



BUNEA, F., A. TSYBAKOV et M. WEGKAMP. "Aggregation for Gaussian regression". In : *The Annals of Statistics* 35.4 (2007), p. 1674–1697.



CESA-BIANCHI, N. et G. LUGOSI. *Prediction, Learning, and Games*. Cambridge University Press, 2006.



GAILLARD, P. et 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), p. 729–769.



LECUÉ, G. "Optimal Rates of Aggregation in Classification Under Low Noise Assumption". Thèse de doct. 2006.



PIERRE, A., C. VINCENT et L. GUILLAUME. "Estimation bounds and sharp oracle inequalities of regularized procedures with Lipschitz loss functions". In : *arXiv preprint arXiv :1702.01402* (2017).



RIGOLLET, P. et A. TSYBAKOV. "Exponential screening and optimal rates of sparse estimation". In : *The Annals of Statistics* (2011), p. 731–771.



STEINHARDT, J., S. WAGER et P. LIANG. "The Statistics of Streaming Sparse Regression". In : *arXiv preprint arXiv :1412.4182* (2014).



XIAO, L. "Dual averaging methods for regularized stochastic learning and online optimization". In : *Journal of Machine Learning Research* 11 (2010), p. 2543–2596.