# a Randomized Proximal Gradient Method with Structure-Adapted Sampling

Franck Iutzeler LJK, Univ. Grenoble Alpes

Journées SMAI MODE 2020 7-9 Sept. 2020



Structure	Regularization
sparsity	$r = \  \cdot \ _1$
anti-sparsity	$r = \  \cdot \ _{\infty}$
low rank	$r = \  \cdot \ _{*}$
•	•
•	•
•	•

Linear inverse problems: for a chosen regularization, we seek

 $x^* \in \arg\min_x r(x)$  such that Ax = b

#### Regularized Empirical Risk Minimization problem:

Find	$x^\star \in rg\min_{x \in \mathbb{R}^n}$	$\mathcal{R}\left(x; \{a_i, b_i\}_{i=1}^m\right)$	+	$\lambda r(x)$
		obtained from		chosen
		statistical modeling		regularization

e.g. Lasso: Find  $x^{\star} \in rgmin_{x \in \mathbb{R}^n} \quad \sum_{i=1}^m \frac{1}{2} (a_i^{\top} x - b_i)^2 \quad + \quad \lambda \; \|x\|_1$ 

Regularization can improve statistical properties (generalization, stability, ...).

♦ Tibshirani: Regression shrinkage and selection via the lasso. Journal of the Royal Statistical Society (1996)

♦ Tibshirani *et al.*: Sparsity and smoothness via the fused lasso. Journal of the Royal Statistical Society (2004)

Vaiter, Peyré, Fadili: Model consistency of partly smooth regularizers. IEEE Trans. on Information Theory (2017)

#### **Composite minimization**

Find
$$x^* \in \arg \min_{x \in \mathbb{R}^n}$$
 $\mathcal{R}(x; \{a_i, b_i\}_{i=1}^m)$  $+ \lambda r(x)$ Find $x^* \in \arg \min_{x \in \mathbb{R}^n}$  $f(x) + g(x)$ smoothnon-smooth

- > f: differentiable surrogate of the empirical risk ⇒ Gradient non-linear smooth function that depends on all the data
- > g: non-smooth but chosen regularization ⇒ Proximity operator non-differentiability on some manifolds implies structure on the solutions

closed form/easy for many regularizations:

$$\mathbf{prox}_{\gamma g}(u) = \arg\min_{y \in \mathbb{R}^n} \left\{ g(y) + \frac{1}{2\gamma} \|y - u\|_2^2 \right\} \qquad \begin{array}{l} -g(x) = \|x\|_1 \\ -g(x) = TV(x) \\ -g(x) = indicator_C(x) \end{array}$$

Natural optimization method: proximal gradient

$$\begin{cases} u_{k+1} = x_k - \gamma \nabla f(x_k) \\ x_{k+1} = \mathbf{prox}_{\gamma g}(u_{k+1}) \end{cases}$$

and its stochastic variants: proximal sgd, etc.

#### Example: LASSO

Find
$$x^* \in \arg\min_{x \in \mathbb{R}^n}$$
 $\mathcal{R}(x; \{a_i, b_i\}_{i=1}^m)$  $+ \lambda r(x)$ Find $x^* \in \arg\min_{x \in \mathbb{R}^n}$  $\frac{1}{2} ||Ax - b||_2^2$  $+ \lambda ||x||_1$ smoothnon-smooth

Coordinates	Structure	$\leftrightarrow$	Optimality conditions
$\forall i$	$x_i^{\star} = 0$	$\Leftrightarrow$	$A_i^{\top}(Ax^{\star} - b) \in [-\lambda, \lambda]$

 $\begin{array}{l} \mbox{Proximity Operator: per coordinate} \\ \left[ \mbox{prox}_{\gamma\lambda \|\cdot\|_1}(u) \right]_i = \left\{ \begin{array}{l} u_i - \lambda\gamma & \mbox{if } u_i > \lambda\gamma \\ 0 & \mbox{if } u_i \in [-\lambda\gamma;\lambda\gamma] \\ u_i + \lambda\gamma & \mbox{if } u_i < -\lambda\gamma \end{array} \right. \end{array}$ 

Proximal Gradient (aka ISTA):

$$\begin{cases} u_{k+1} = x_k - \gamma A^\top (Ax_k - b) \\ x_{k+1} = \mathbf{prox}_{\gamma \lambda \parallel \cdot \parallel_1} (u_{k+1}) \end{cases}$$



#### **Example: LASSO**

Find 
$$x^* \in \arg\min_{x \in \mathbb{R}^n} \mathcal{R}(x; \{a_i, b_i\}_{i=1}^m) + \lambda r(x)$$
  
Find  $x^* \in \arg\min_{x \in \mathbb{R}^n} \frac{1}{2} ||Ax - b||_2^2 + \lambda ||x||_1$   
smooth non-smooth

 $\begin{array}{rcl} \text{Coordinates} & \textbf{Structure} & \leftrightarrow & \textbf{Optimality conditions} & \leftrightarrow & \textbf{Proximity operation} \\ \forall i & x_i^* = \mathbf{0} & \Leftrightarrow & A_i^\top (Ax^* - b) \in [-\lambda, \lambda] & \Leftrightarrow & \begin{bmatrix} \textbf{prox}_{\gamma \lambda \| \cdot \|_1}(u^*) \end{bmatrix}_i = \mathbf{0} \\ & u^* = x^* - \gamma A^\top (Ax^* - b) \end{array}$ 



Proximal Gradient (aka ISTA):

$$\begin{pmatrix} u_{k+1} = x_k - \gamma A^\top (Ax_k - b) \\ x_{k+1} = \mathbf{prox}_{\gamma \lambda \parallel \cdot \parallel_1} (u_{k+1}) \end{pmatrix}$$



Iterates  $(x_k)$  reach the same structure as  $x^*$  in finite time!

#### >>> Mathematical properties of Proximal Algorithms



> project on manifolds

Let  $\mathcal{M}$  be a manifold and  $u_k$  such that

$$x_k = \mathbf{prox}_{\gamma g}(u_k) \in \mathcal{M}$$
 and  $\frac{u_k - x_k}{\gamma} \in \operatorname{ri} \partial g(x_k)$ 

If *g* is partly smooth at  $x_k$  relative to  $\mathcal{M}$ , then

 $\mathbf{prox}_{\gamma g}(u) \in \mathcal{M}$ 

#### for any u close to $u_k$ .

- Hare, Lewis: Identifying active constraints via partial smoothness and prox-regularity. Journal of Convex Analysis (2004)
- Daniilidis, Hare, Malick: Geometrical interpretation of the predictor-corrector type algorithms in structured optimization problems. Optimization (2006)

#### >>> Mathematical properties of Proximal Algorithms



> project on manifolds> identify the optimal structure

Let  $(x_k)$  and  $(u_k)$  be a pair of sequences such that

$$x_k = \mathbf{prox}_{\gamma g}(u^k) \to x^\star = \mathbf{prox}_{\gamma g}(u^\star)$$

and  $\mathcal{M}$  be a manifold. If  $x^{\star} \in \mathcal{M}$  and

 $\exists \varepsilon > 0 \text{ such that for all } u \in \mathcal{B}(u^*, \varepsilon), \ \mathbf{prox}_{\gamma g}(u) \in \mathcal{M}$  (QC)

holds, then, after some finite but unknown time,  $x_k \in \mathcal{M}$ .

- ◊ Lewis: Active sets, nonsmoothness, and sensitivity. SIAM Journal on Optimization (2002)
- Fadili, Malick, Peyré: Sensitivity analysis for mirror-stratifiable convex functions. SIAM Journal on Optimization (2018)
- Bareilles, I.: On the Interplay between Acceleration and Identification for the Proximal Gradient algorithm. Computational Optimization and Applications (2020)

#### >>> "Nonsmoothness can help"

#### > Nonsmoothness is actively studied in Numerical Optimization...

Subgradients, Partial Smoothness/prox-regularity, Bregman metrics, Error Bounds/Kurdyka-Łojasiewicz, etc.

- Hare, Lewis: Identifying active constraints via partial smoothness and prox-regularity. Journal of Convex Analysis (2004)
- ◇ Lemarechal, Oustry, Sagastizabal: The U-Lagrangian of a convex function. Transactions of the AMS (2000)
- Bolte, Daniilidis, Lewis: The Łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems. SIAM Journal on Optimization (2007)
- Chen, Teboulle: A proximal-based decomposition method for convex minimization problems. Mathematical Programming (1994)

> Nonsmoothness is actively studied in Numerical Optimization...

Subgradients, Partial Smoothness/prox-regularity, Bregman metrics, Error Bounds/Kurdyka-Łojasiewicz, etc.

> ...but often suffered because of lack of structure/expression. Bundle methods, Gradient Sampling, Smoothing, Inexact proximal methods, etc.

- ◊ Nesterov: Smooth minimization of non-smooth functions. Mathematical Programming (2005)
- Burke, Lewis, Overton: A robust gradient sampling algorithm for nonsmooth, nonconvex optimization. SIAM Journal on Optimization (2005)
- ◊ Solodov, Svaiter: A hybrid projection-proximal point algorithm. Journal of convex analysis (1999)
- ♦ de Oliveira, Sagastizábal: Bundle methods in the XXIst century: A bird's-eye view. Pesquisa Operacional (2014)

- > Nonsmoothness is actively studied in Numerical Optimization... Subgradients, Partial Smoothness/prox-regularity, Bregman metrics, Error Bounds/Kurdvka-Łoiasiewicz. etc.
- > ...but often suffered because of lack of structure/expression. Bundle methods, Gradient Sampling, Smoothing, Inexact proximal methods, etc.
- > For Machine Learning objectives, it can often be harnessed
  - Explicit/"proximable" regularizations l1, nuclear norm
  - We know the expressions and activity of sought structures sparsity, rank

- ♦ Bach, et al.: Optimization with sparsity-inducing penalties. Foundations and Trends in Machine Learning (2012)
- ◊ Massias, Salmon, Gramfort: Celer: a fast solver for the lasso with dual extrapolation. ICML (2018)
- Liang, Fadili, Peyré: Local linear convergence of forward-backward under partial smoothness. NeurIPS (2014)
- O'Donoghue, Candes: Adaptive restart for accelerated gradient schemes. Foundations of computational mathematic (2015)

Find
$$x^* \in \arg \min_{x \in \mathbb{R}^n}$$
 $\mathcal{R}(x; \{a_i, b_i\}_{i=1}^m)$  $+ \lambda r(x)$ Find $x^* \in \arg \min_{x \in \mathbb{R}^n}$  $f(x)$  $+ g(x)$ smoothnon-smooth

A reason why the nonsmoothness of ML problems can be leveraged is that we <u>know</u> the *expressions* and *activity* of sought structures.

Mathematically, we can design a *lookout collection*  $C = \{M_1, ..., M_p\}$  of closed sets such that:

- (i) we have a projection mapping  $\text{proj}_{\mathcal{M}_i}$  onto  $\mathcal{M}_i$  for all *i*;
- (ii) **prox**<sub> $\gamma g$ </sub>(*u*) is a singleton and can be computed explicitly for any *u* and  $\gamma$ ;

(iii) upon computation of  $x = \mathbf{prox}_{\gamma g}(u)$ , we know if  $x \in \mathcal{M}_i$  or not for all *i*.

 $\Rightarrow$  Structure can be directly observed.

Example: Sparse structure and  $g = \|\cdot\|_1, \|\cdot\|_{0.5}^{0.5}, \|\cdot\|_0, \dots$ 

$$\mathsf{C} = \{\mathcal{M}_1, \dots, \mathcal{M}_n\} \quad ext{ with } \mathcal{M}_i = \{x \in \mathbb{R}^n : x_i = 0\}$$

> Identification of proximal algorithms

 $\Rightarrow$  After some finite time, the iterates reach the optimal structure.

> Problems with noticeable strucure

 $\Rightarrow$  Structure can be directly observed.

Let's use the structure progressively uncovered in proximal gradient!

not necessarily close to  $x^*$ , no explicit or bad dual, when screening methods are difficult to implement

> Identification of proximal algorithms

 $\Rightarrow$  After some finite time, the iterates reach the optimal structure.

> Problems with noticeable strucure
 ⇒ Structure can be directly observed.

Let's use the structure progressively uncovered in proximal gradient! not necessarily close to x<sup>\*</sup>, no explicit or bad dual, when screening methods are difficult to implement

**Idea:** Define  $\mathcal{M}_k = \mathbb{R}^n \bigcap_{i:x_k \in \mathcal{M}_i} \mathcal{M}_i$  and  $\mathcal{M}^* := \mathbb{R}^n \bigcap_{i:x^* \in \mathcal{M}_i} \mathcal{M}_i$ , then:

 $\mathcal{M}_k \subset \mathbb{R}^n$  partial identif/screening and  $\mathcal{M}_k = \mathcal{M}^\star$  after some finite time identification

Can we reduce the dimension of the problem on the fly using  $M_k$ ?

## Proximal Gradient methods with Adaptive Subspace Sampling

Dmitry Grishchenko F. I.



PhD student Univ. Grenoble Alpes Jerôme Malick



CNRS and LJK

to appear in Mathematics of Operations Research

https://arxiv.org/abs/2004.13356

<u>Disclaimer</u>: We assume that the identified manifolds are linear subspaces  $eg: \|Dx\|_1$  but not separability of *g*.

$$egin{array}{rcl} y_k &=& x_k - \gamma 
abla f(x_k) \ z_k &=& y_k \ x_{k+1} &= \mathbf{prox}_{\gamma g}(z_k) \end{array}$$

> Vanilla Proximal gradient identifies but does not use it full gradient computed at each iteration

Example: Sparse structure and  $g = \| \cdot \|_1$ 

$$C = \{\mathcal{M}_1, \dots, \mathcal{M}_n\}$$
 with  $\mathcal{M}_i = \{x \in \mathbb{R}^n : x_i = 0\}$ 

<u>Disclaimer</u>: We assume that the identified manifolds are linear subspaces  $eg: ||Dx||_1$  but not separability of *g*.

Observe 
$$\mathcal{M}_k = \mathbb{R}^n \bigcap_{i:x_k \in \mathcal{M}_i} \mathcal{M}_i$$
  
 $y_k = x_k - \gamma \nabla f(x_k)$   
 $z_k = \operatorname{proj}_{\mathcal{M}_k}(y_k) + \operatorname{proj}_{\mathcal{M}_k}^{\perp}(z_{k-1})$   
 $x_{k+1} = \operatorname{prox}_{\gamma g}(z_k)$ 

> Direct Use of Identification may not converge eg: starting with 0

Example: Sparse structure and  $g = \| \cdot \|_1$ 

$$C = \{\mathcal{M}_1, \dots, \mathcal{M}_n\} \quad \text{with } \mathcal{M}_i = \{x \in \mathbb{R}^n : x_i = 0\}$$
$$\mathcal{M}_k = \{x \in \mathbb{R}^n : x_i = 0 \text{ if } x_{i,k} = 0\}$$

<u>Disclaimer</u>: We assume that the identified manifolds are linear subspaces  $eg: ||Dx||_1$  but not separability of *g*.

Observe 
$$\mathcal{M}_k = \mathbb{R}^n \bigcap_{i:x_k \in \mathcal{M}_i} (\xi_{k,i} \mathcal{M}_i + (1 - \xi_{k,i}) \mathbb{R}^n)$$
 for  $\xi_{k,i} \sim \mathcal{B}(p)$   
 $y_k = x_k - \gamma \nabla f(x_k)$   
 $z_k = \operatorname{proj}_{\mathcal{M}_k}(y_k) + \operatorname{proj}_{\mathcal{M}_k}^{\perp}(z_{k-1})$   
 $x_{k+1} = \operatorname{prox}_{\gamma g}(z_k)$ 

> Mixing Identification and Randomized "coordinate" descent biases convergence issues

Example: Sparse structure and  $g = \| \cdot \|_1$ 

$$C = \{\mathcal{M}_1, \dots, \mathcal{M}_n\} \text{ with } \mathcal{M}_i = \{x \in \mathbb{R}^n : x_i = 0\}$$
$$\mathcal{M}_k = \{x \in \mathbb{R}^n : x_i = 0 \text{ if } x_{i,k} = 0 \text{ for some } i\}$$

#### >>> Adaptive Proximal Gradient

<u>Disclaimer</u>: We assume that the identified manifolds are linear subspaces  $eg: ||Dx||_1$  but not separability of *g*.

Observe 
$$\mathcal{M}_k = \mathbb{R}^n \bigcap_{i:x_k \in \mathcal{M}_i} (\xi_{k,i}\mathcal{M}_i + (1 - \xi_{k,i})\mathbb{R}^n)$$
 for  $\xi_{k,i} \sim \mathcal{B}(p)$   
and compute  $\mathbf{P}_k = \mathbb{E} \operatorname{proj}_{\mathcal{M}_k}$  and  $Q_k = (\mathbf{P}_k)^{-1/2}$   
 $y_k = Q_k (x_k - \gamma \nabla f(x_k))$   
 $z_k = \operatorname{proj}_{\mathcal{M}_k}(y_k) + \operatorname{proj}_{\mathcal{M}_k}^{\perp}(z_{k-1})$   
 $x_{k+1} = \operatorname{prox}_{\gamma g}(Q_k^{-1} z_k)$ 

> Unbiasing with  $Q_k = (\mathbb{E} \text{proj}_{\mathcal{M}_k})^{-1/2}$  works after identification but before... no, which prevents identification...



#### >>> Adaptive Proximal Gradient

<u>Disclaimer</u>: We assume that the identified manifolds are linear subspaces  $eg: ||Dx||_1$  but not separability of *g*.

Observe 
$$\mathcal{M}_k = \mathbb{R}^n \bigcap_{i:x_\ell \in \mathcal{M}_i} (\xi_{k,i}\mathcal{M}_i + (1 - \xi_{k,i})\mathbb{R}^n)$$
 for  $\xi_{k,i} \sim \mathcal{B}(p)$   
and compute  $\mathbf{P}_k = \mathbb{E}\text{proj}_{\mathcal{M}_k}$  and  $Q_k = (\mathbf{P}_k)^{-1/2}$  sometimes, else  $\mathcal{M}_k = \mathcal{M}_{k-1}$   
 $y_k = Q_k (x_k - \gamma \nabla f(x_k))$   
 $z_k = \text{proj}_{\mathcal{M}_k}(y_k) + \text{proj}_{\mathcal{M}_k}^{\perp}(z_{k-1})$   
 $x_{k+1} = \mathbf{prox}_{\gamma g}(Q_k^{-1}z_k)$ 

### > Structure adaptation can be performed at *some* iterations

depends on the amount of change  $\|Q_{k-1}Q_k^{-1}\|$  and harshness of the sparsification  $\lambda_{\min}(Q_k)$ 



$$\begin{aligned} \mathcal{M}_k &= \mathbb{R}^n \bigcap_i (\xi_{k,i} \mathcal{M}_i + (1 - \xi_{k,i}) \mathbb{R}^n) \text{ for } \xi_{k,i} \sim \mathcal{B}(p)(iid) \\ \text{Compute (once) } \mathbf{P} &= \mathbb{E} \text{proj}_{\mathcal{M}_k} \text{ and } Q = (\mathbf{P})^{-1/2} \\ y_k &= Q \left( x_k - \gamma \nabla f(x_k) \right) \\ z_k &= \text{proj}_{\mathcal{M}_k}(y_k) + \text{proj}_{\mathcal{M}_k}^{\perp}(z_{k-1}) \\ z_{k+1} &= \mathbf{prox}_{\gamma g}(Q^{-1}z_k) \end{aligned}$$

**Theorem** Let f be L-smooth and  $\mu$ -strongly convex and g be convex, lsc. Take any  $\gamma \in (0, 2/(\mu + L)]$ . Then, the sequence  $(x_k)$  converges almost surely to the minimizer  $x^*$  of f + g and

$$\mathbb{E}[\|x_k - x^\star\|^2] \leq \left(1 - \lambda_{\min}(\mathbf{P}) rac{2\gamma \mu L}{\mu + L}
ight)^k C$$

$$\begin{aligned} \text{Proof. Let } z^{\star} &= y^{\star} = Q(x^{\star} - \gamma \nabla f(x^{\star})), \\ \mathbb{E}[\|x_{k+1} - x^{\star}\|^{2} |\mathcal{F}_{k-1}] &\leq \mathbb{E}[\|z_{k} - z^{\star}\|^{2} |\mathcal{F}_{k-1}] \\ &\leq \mathbb{E}[\|z_{k-1} - z^{\star}\|^{2} + \mathbb{E}[\|y_{k} - y^{\star}\|^{2}_{\mathbf{P}} |\mathcal{F}_{k-1}] - \|z_{k-1} - z^{\star}\|^{2} \\ &\leq \left(1 - \lambda_{\min}(\mathbf{P}) \frac{2\gamma \mu L}{\mu + L}\right) \|z_{k-1} - z^{\star}\|^{2} \end{aligned}$$

 $\begin{aligned} & \text{if adaptation,} \quad \ell \leftarrow \ell + 1, x^{base} \leftarrow x_{k-c_{\ell}} \\ & \text{(upd. space changes,} \quad \mathcal{M}_{k} = \mathbb{R}^{n} \bigcap_{i:x_{\ell}^{base} = 0} (\xi_{k,i}\mathcal{M}_{i} + (1 - \xi_{k,i})\mathbb{R}^{n}) \text{ for } \xi_{k,i} \sim \mathcal{B}(p)(iid) \\ & \text{but same dist. whithin)} \quad \text{Compute (at each adapt.) } \mathbf{P}_{\ell} = \mathbb{E}\text{proj}_{\mathcal{M}_{\ell}} \text{ and } Q_{\ell} = (\mathbf{P}_{\ell})^{-1/2} \\ & y_{k} \quad = Q_{\ell} \left( x_{k} - \gamma \nabla f(x_{k}) \right) \\ & z_{k} \quad = \text{proj}_{\mathcal{M}_{k}}(y_{k}) + \text{proj}_{\mathcal{M}_{k}}^{\perp}(z_{k-1}) \\ & x_{k+1} \quad = \mathbf{prox}_{\gamma g}(Q_{\ell}^{-1}z_{k}) \end{aligned}$ 

**Theorem** Let f be L-smooth and  $\mu$ -strongly convex and g be convex, lsc. Take any  $\gamma \in (0, 2/(\mu + L)]$ . Consider the following adaptation strategy: 1) If the structure of  $x_k$  changes, choose a new sampling  $x_k \to x^{base}$ ,  $\lambda_{\min}(\mathbf{P}_{\ell}) \ge \lambda$ 

2) Compute  $c_{\ell}$  so that  $\|Q_{\ell}Q_{\ell-1}^{-1}\|^2 \left(1 - \lambda_{\min}(\mathbf{P}_{\ell-1})\frac{2\gamma\mu L}{\mu+L}\right)^{c_{\ell}} \leq \left(1 - \lambda\frac{\gamma\mu L}{\mu+L}\right)$ 

3) Apply the new sampling after  $c_{\ell}$  iterations

Then, the sequence  $(x_k)$  converges almost surely to the minimizer  $x^*$  of f + g and

$$\mathbb{E}[\|x_k - x^\star\|^2] \le \left(1 - \lambda \frac{\gamma \mu L}{\mu + L}\right)^\ell C$$

where  $\ell$  is the number of adaptations before k.

if adaptation,  $\ell \leftarrow \ell + 1, x^{base} \leftarrow x_{k-c_{\ell}}$ (upd. space changes,  $\mathcal{M}_{k} = \mathbb{R}^{n} \bigcap_{i:x_{i}^{base}=0} (\xi_{k,i}\mathcal{M}_{i} + (1 - \xi_{k,i})\mathbb{R}^{n})$  for  $\xi_{k,i} \sim \mathcal{B}(p)(iid)$ but same dist. whithin) Compute (at each adapt.)  $\mathbf{P}_{\ell} = \mathbb{E}\mathrm{proj}_{\mathcal{M}_{\ell}}$  and  $Q_{\ell} = (\mathbf{P}_{\ell})^{-1/2}$   $y_{k} = Q_{\ell} (x_{k} - \gamma \nabla f(x_{k}))$   $z_{k} = \mathrm{proj}_{\mathcal{M}_{k}}(y_{k}) + \mathrm{proj}_{\mathcal{M}_{k}}^{\perp}(z_{k-1})$  $x_{k+1} = \mathbf{prox}_{\gamma g}(Q_{\ell}^{-1}z_{k})$ 

1) If the structure of  $x_k$  changes, choose a new sampling  $x_k \to x^{base}$ ,  $\lambda_{\min}(\mathbf{P}_{\ell}) \ge \lambda$ 2) Compute  $c_{\ell}$  so that  $\|Q_{\ell}Q_{\ell-1}^{-1}\|^2 \left(1 - \lambda_{\min}(\mathbf{P}_{\ell-1})\frac{2\gamma\mu L}{\mu+L}\right)^{c_{\ell}} \le \left(1 - \lambda\frac{\gamma\mu L}{\mu+L}\right)$ 3) Apply the new sampling after  $c_{\ell}$  iterations

We "compensate" as along as the structure changes

**Theorem** Under the same assumptions as before + (QC), the rate improves to

$$\mathbb{E}[\|\boldsymbol{x}_k - \boldsymbol{x}^\star\|^2] \leq \left(1 - \lambda_{\min}(\mathbf{P}^\star) \frac{2\gamma \mu L}{\mu + L}\right)^k C$$

linear in *iterations* (not adapt.) with a rate depending on the final structure.

### TV-reg. logistic regression on a1a (1605 $\times$ 143), 90% final *jump* sparsity



- > Iterate structure enforced by nonsmooth regularizers can be used to adapt the selection probabilities of coordinate descent/sketching;
- > Before identification, adaptation *has to be moderate*;
- > The cost of adaption may be big (in iterations and computation) but this can be mitigated in many practical cases.

- > Machine Learning problems often have a *noticeable structure*;
- > We can design a *lookout collection* C = {*M*<sub>1</sub>, ..., *M<sub>p</sub>*} of sets: (i) with easy projections; (ii) identified by proximity operations; (iii) we *know* if these sets are identified or not;
- > This structure can/should be harnessed but may be tricky before identification.

▷ Malick & I.: Nonsmoothness in Machine Learning: specific structure, proximal identification, and applications, review/pedagogical paper coming hopefully soon



Thank you! - Franck IUTZELER http://www.iutzeler.org