

# Proximal Identification and Applications

Jérôme MALICK

CNRS, Lab. J. Kuntzmann, Grenoble (France)



Workshop Optimization for Machine Learning – Luminy – March 2020

talk based on materiel from joint work with

G. Peyré



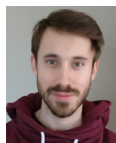
J. Fadili



G. Garrigos



F. Iutzeler



D. Grishchenko

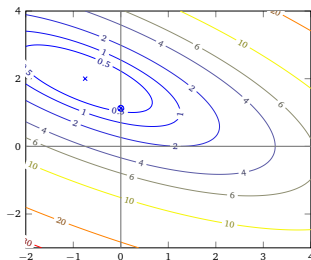


## Example of stability

$$\min_{x \in \mathbb{R}^d} \frac{1}{2} \|Ax - y\|^2 + \lambda \|x\|_1 \quad (\text{LASSO})$$

**Stability:** the support of optimal solutions is stable under small perturbations

**Illustration** (on an instance with  $d = 2$ )

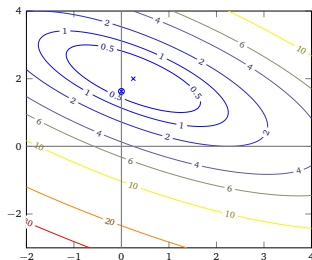
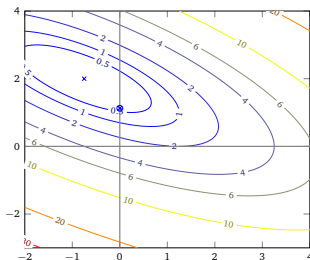


## Example of stability

$$\min_{x \in \mathbb{R}^d} \frac{1}{2} \|Ax - y\|^2 + \lambda \|x\|_1 \quad (\text{LASSO})$$

**Stability:** the support of optimal solutions is stable under small perturbations

**Illustration** (on an instance with  $d = 2$ )

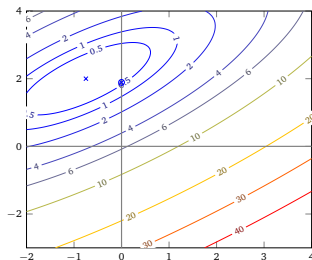
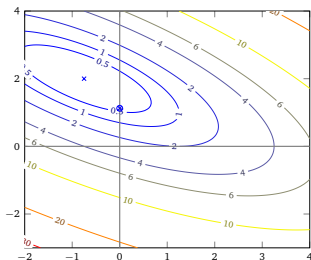


## Example of stability

$$\min_{x \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{A}x - y\|^2 + \lambda \|x\|_1 \quad (\text{LASSO})$$

**Stability:** the support of optimal solutions is stable under small perturbations

**Illustration** (on an instance with  $d = 2$ )

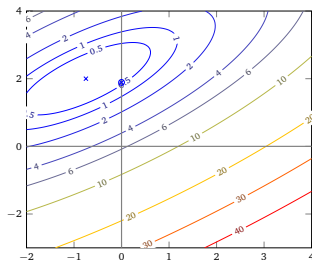
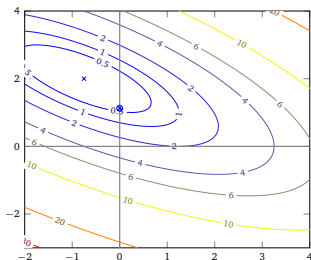


## Example of stability

$$\min_{x \in \mathbb{R}^d} \quad \frac{1}{2} \|Ax - y\|^2 + \lambda \|x\|_1 \quad (\text{LASSO})$$

**Stability:** the support of optimal solutions is stable under small perturbations

Illustration (on an instance with  $d = 2$ )



More generally: [Lewis '02] sensitivity analysis of partly-smooth functions

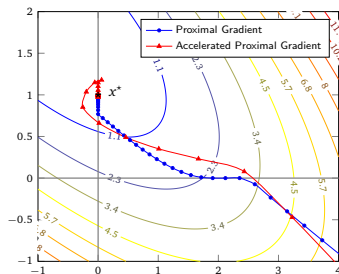
(remind Clarice's talk, this morning)

## Example of identification

$$\min_{x \in \mathbb{R}^d} \frac{1}{2} \|Ax - y\|^2 + \lambda \|x\|_1 \quad (\text{LASSO})$$

**Identification:** (proximal-gradient) algorithms produce iterates...

...that eventually have the same support as the optimal solution



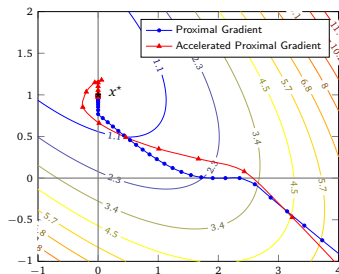
Runs of two proximal-gradient algos  
(same instance with  $d = 2$ )

## Example of identification

$$\min_{x \in \mathbb{R}^d} \frac{1}{2} \|Ax - y\|^2 + \lambda \|x\|_1 \quad (\text{LASSO})$$

**Identification:** (proximal-gradient) algorithms produce iterates...

...that eventually have the same support as the optimal solution



Runs of two proximal-gradient algos  
(same instance with  $d = 2$ )

Well-studied, see e.g. [Bertsekas '76], [Wright '96], [Lewis Drusvyatskiy '13]...

# Outline

- 1 General stability of regularized problems
- 2 Enlarged identification of proximal algorithms
- 3 Application: communication-efficient federated learning
- 4 Application: model consistency for regularized least-squares

# Outline

- 1 General stability of regularized problems
- 2 Enlarged identification of proximal algorithms
- 3 Application: communication-efficient federated learning
- 4 Application: model consistency for regularized least-squares

## Stability or sensitivity analysis

Parameterized composite optimization problem (smooth + nonsmooth)

$$\min_{x \in \mathbb{R}^d} F(x, p) + R(x),$$

Typically nonsmooth  $R$  traps solutions in low-dimensional manifolds

**Stability:** Optimal solutions lie on a manifold:  $x^*(p) \in M$  for  $p \sim p_0$

Studied in e.g. [Hare Lewis '10] [Vaier *et al* '15] [Liang *et al* '16]...

Example 1:  $R = \|\cdot\|_1$ ,  $\text{supp}(x^*(p)) = \text{supp}(x^*(p_0))$

# Stability or sensitivity analysis

Parameterized composite optimization problem (smooth + nonsmooth)

$$\min_{x \in \mathbb{R}^d} F(x, p) + R(x),$$

Typically nonsmooth  $R$  traps solutions in low-dimensional manifolds

**Stability:** Optimal solutions lie on a manifold:  $x^*(p) \in M$  for  $p \sim p_0$

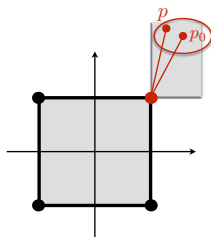
Studied in e.g. [Hare Lewis '10] [Vaier *et al* '15] [Liang *et al* '16]...

Example 1:  $R = \|\cdot\|_1$ ,  $\text{supp}(x^*(p)) = \text{supp}(x^*(p_0))$

Example 2:  $R = \iota_{\mathbb{B}_\infty}$  (indicator function)

projection onto the  $\ell_\infty$  ball

Many examples in machine learning...

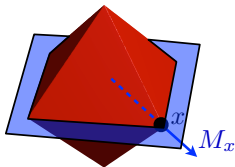


## Structure of nonsmooth regularizers

Many of the regularizers used in machine learning or image processing have a strong primal-dual structure (“mirror-stratifiable” [Fadili, M., Peyré '18]) ...that can be exploit to get (enlarged) stability/identification results

Examples: (associated unit ball and low-dimensional manifold where  $x$  belongs)

- $R = \|\cdot\|_1$  ( and  $\|\cdot\|_\infty$  or other polyedral gauges)



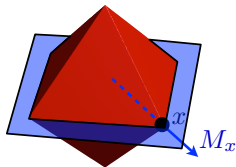
$$M_x = \{z : \text{supp}(z) = \text{supp}(x)\}$$

## Structure of nonsmooth regularizers

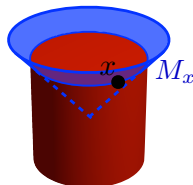
Many of the regularizers used in machine learning or image processing have a strong primal-dual structure (“mirror-stratifiable” [Fadili, M., Peyré '18]) ...that can be exploited to get (enlarged) stability/identification results

Examples: (associated unit ball and low-dimensional manifold where  $x$  belongs)

- $R = \|\cdot\|_1$  ( and  $\|\cdot\|_\infty$  or other polyedral gauges)
- nuclear norm (aka trace-norm)  $R(X) = \sum_i |\sigma_i(X)| = \|\sigma(X)\|_1$



$$M_x = \{z : \text{supp}(z) = \text{supp}(x)\}$$



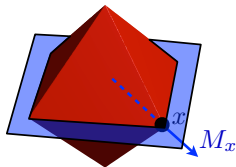
$$M_x = \{z : \text{rank}(z) = \text{rank}(x)\}$$

## Structure of nonsmooth regularizers

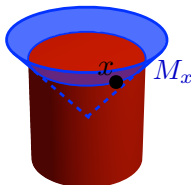
Many of the regularizers used in machine learning or image processing have a strong primal-dual structure (“mirror-stratifiable” [Fadili, M., Peyré '18]) ...that can be exploited to get (enlarged) stability/identification results

Examples: (associated unit ball and low-dimensional manifold where  $x$  belongs)

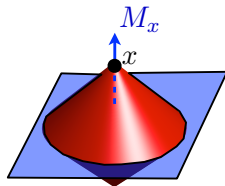
- $R = \|\cdot\|_1$  ( and  $\|\cdot\|_\infty$  or other polyedral gauges)
- nuclear norm (aka trace-norm)  $R(X) = \sum_i |\sigma_i(X)| = \|\sigma(X)\|_1$
- group- $\ell_1$   $R(x) = \sum_{b \in \mathcal{B}} \|x_b\|_2$  ( e.g.  $R(x) = \|x_{1,2}\| + |x_3|$  )



$$M_x = \{z : \text{supp}(z) = \text{supp}(x)\}$$



$$M_x = \{z : \text{rank}(z) = \text{rank}(x)\}$$



$$M_x = \{0\} \times \{0\} \times \mathbb{R}$$

## Recall on stratifications

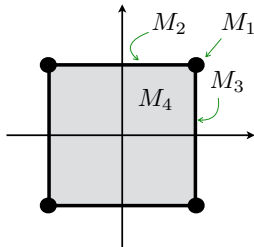
A stratification of a set  $D \subset \mathbb{R}^d$  is a (finite) partition  $\mathcal{M} = \{M_i\}_{i \in I}$

$$D = \bigcup_{i \in I} M_i$$

with so-called “strata” (e.g. smooth/affine manifolds) which fit nicely:

$$M \cap \text{cl}(M') \neq \emptyset \implies M \subset \text{cl}(M')$$

**Example:**  $\mathbb{B}_\infty$  the unit  $\ell_\infty$ -ball in  $\mathbb{R}^2$   
a stratification with 9 (affine) strata



Other examples: “tame” sets, remind Edouard’s talk

## Recall on stratifications

A stratification of a set  $D \subset \mathbb{R}^d$  is a (finite) partition  $\mathcal{M} = \{M_i\}_{i \in I}$

$$D = \bigcup_{i \in I} M_i$$

with so-called “strata” (e.g. smooth/affine manifolds) which fit nicely:

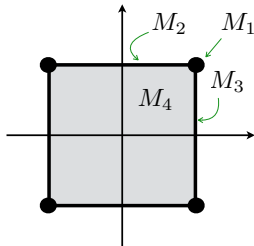
$$M \cap \text{cl}(M') \neq \emptyset \implies M \subset \text{cl}(M')$$

This relation induces a (partial) ordering  $M \leq M'$

**Example:**  $\mathbb{B}_\infty$  the unit  $\ell_\infty$ -ball in  $\mathbb{R}^2$   
a stratification with 9 (affine) strata

$$M_1 \leq M_2 \leq M_4$$

$$M_1 \leq M_3 \leq M_4$$



Other examples: “tame” sets, remind Edouard’s talk

## Mirror-stratifiable regularizations

(primal) stratification  $\mathcal{M} = \{M_i\}_{i \in I}$  and (dual) stratification  $\mathcal{M}^* = \{M_i^*\}_{i \in I}$

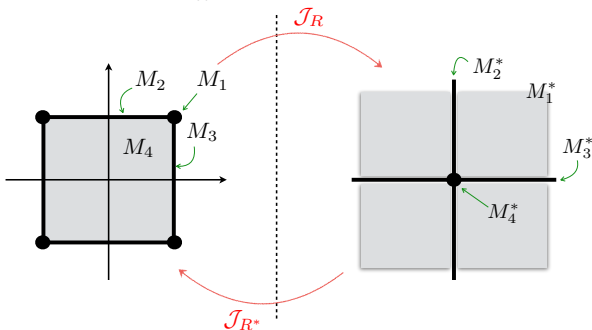
in **one-to-one decreasing correspondence**

through the transfert operator  $\mathcal{J}_R(S) = \bigcup_{x \in S} \text{ri}(\partial R(x))$

Simple example:

$$R = \iota_{\mathbb{B}_\infty}$$

$$R^* = \|\cdot\|_1$$



$$\mathcal{J}_R(M_i) = \bigcup_{x \in M_i} \text{ri} \partial R(x) = \text{ri} N_{\mathbb{B}_\infty}(x) = M_i^* \quad M_i = \text{ri} \partial \|x\|_1 = \bigcup_{x \in M_i^*} \text{ri} \partial R^*(x) = \mathcal{J}_{R^*}(M_i^*)$$

# Enlarged stability result

Theorem (Fadili, M., Peyré '18)

For the composite optimization problem (smooth + nonsmooth)

$$\min_{x \in \mathbb{R}^d} F(x, p) + R(x),$$

satisfying mild assumptions (unique minimizer  $x^*(p_0)$  at  $p_0$  and objective uniformly level-bounded in  $x$ ), if  $R$  is mirror-stratifiable, then for  $p \sim p_0$ ,

$$M_{x^*(p_0)} \leq M_{x^*(p)} \leq \mathcal{J}_{R^*}(M_{u^*(p_0)}^*)$$

If  $R = \|\cdot\|_1$ , then  $\text{supp}(x^*(p_0)) \subseteq \text{supp}(x^*(p)) \subseteq \{i : |u^*(p_0)_i| = 1\}$

## Enlarged stability result

Theorem (Fadili, M., Peyré '18)

For the composite optimization problem (smooth + nonsmooth)

$$\min_{x \in \mathbb{R}^d} F(x, p) + R(x),$$

satisfying mild assumptions (unique minimizer  $x^*(p_0)$  at  $p_0$  and objective uniformly level-bounded in  $x$ ), if  $R$  is mirror-stratifiable, then for  $p \sim p_0$ ,

$$M_{x^*(p_0)} \leq M_{x^*(p)} \leq \mathcal{J}_{R^*}(M_{u^*(p_0)}^*)$$

If  $R = \|\cdot\|_1$ , then  $\text{supp}(x^*(p_0)) \subseteq \text{supp}(x^*(p)) \subseteq \{i : |u^*(p_0)_i| = 1\}$

Remark: Optimality conditions for a primal-dual solution  $(x^*(p), u^*(p))$

$$u^*(p) = -\nabla F(x^*(p), p) \in \partial R(x^*(p))$$

In the non-degenerate case:  $u^*(p_0) \in \text{ri}(\partial R(x^*(p_0)))$

$$M_{x^*(p_0)} = M_{x^*(p)} \quad (= \mathcal{J}_{R^*}(M_{u^*(p_0)}^*))$$

we have the exact stability, expected [Lewis '02]

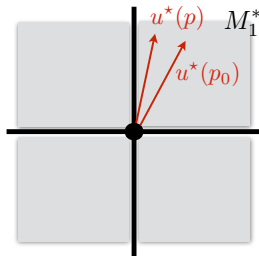
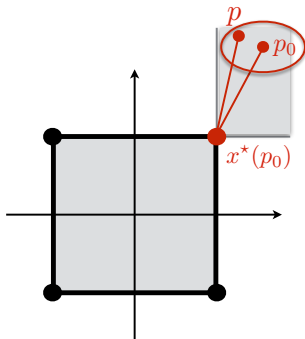
## Enlarged stability illustrated

$$\begin{cases} \min & \frac{1}{2} \|x - p\|^2 \\ & \|x\|_\infty \leq 1 \end{cases}$$

$$\begin{cases} \min & \frac{1}{2} \|u - p\|^2 + \|u\|_1 \\ & u \in \mathbb{R}^n \end{cases}$$

**Non-degenerate case:**  $u^*(p_0) = p_0 - x^*(p_0) \in \text{ri } N_{\mathbb{B}_\infty}(x^*(p_0))$

$$\implies M_1 = M_{x^*(p_0)} = M_{x^*(p)} \quad (\text{in this case } x^*(p) = x^*(p_0))$$

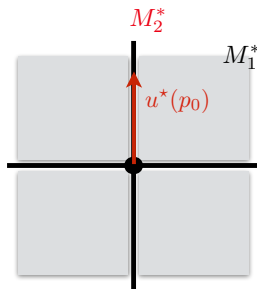
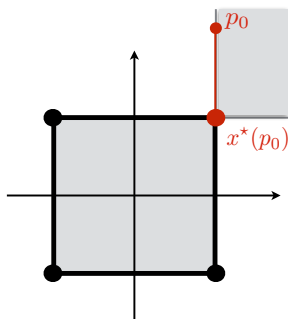


## Enlarged stability illustrated

$$\begin{cases} \min & \frac{1}{2} \|x - p\|^2 \\ & \|x\|_\infty \leq 1 \end{cases}$$

$$\begin{cases} \min & \frac{1}{2} \|u - p\|^2 + \|u\|_1 \\ & u \in \mathbb{R}^n \end{cases}$$

General case:  $u^*(p_0) = p_0 - x^*(p_0) \in \textcolor{red}{\nexists} N_{\mathbb{B}_\infty}(x^*(p))$



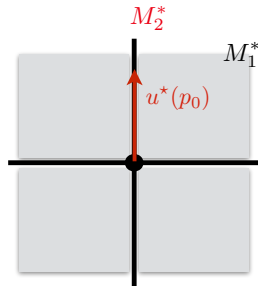
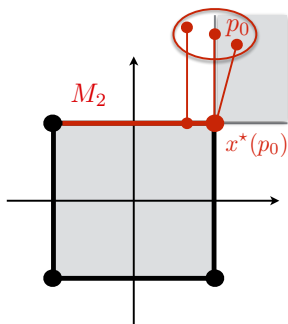
## Enlarged stability illustrated

$$\begin{cases} \min & \frac{1}{2} \|x - p\|^2 \\ & \|x\|_\infty \leq 1 \end{cases}$$

$$\begin{cases} \min & \frac{1}{2} \|u - p\|^2 + \|u\|_1 \\ & u \in \mathbb{R}^n \end{cases}$$

General case:  $u^*(p_0) = p_0 - x^*(p_0) \in \textcolor{red}{\nexists} N_{\mathbb{B}_\infty}(x^*(p))$

$$\implies M_1 = M_{\textcolor{red}{x}^*(p_0)} \leq M_{\textcolor{blue}{x}^*(p)} \leq \mathcal{J}_{R^*}(M_{\textcolor{red}{u}^*(p_0)}^*) = M_2$$



# Outline

- 1 General stability of regularized problems
- 2 Enlarged identification of proximal algorithms
- 3 Application: communication-efficient federated learning
- 4 Application: model consistency for regularized least-squares

## Activity identification

Composite optimization problem (smooth + nonsmooth)

$$\min_{x \in \mathbb{R}^d} F(x) + R(x)$$

Basic proximal-gradient algorithm

$$x_{k+1} = \text{prox}_{\gamma R}(x_k - \gamma \nabla F(x_k))$$

$$\text{prox}_{\gamma R}(x) = \underset{y}{\operatorname{argmin}} R(y) + \frac{1}{2\gamma} \|y - x\|^2$$

$\text{prox}_{\gamma R}(x)$  easy to compute in some important cases

e.g. explicit expression for  $R = \|\cdot\|_1$  (soft-thresholding)

## Activity identification

Composite optimization problem (smooth + nonsmooth)

$$\min_{x \in \mathbb{R}^d} F(x) + R(x)$$

Basic proximal-gradient algorithm

$$x_{k+1} = \text{prox}_{\gamma R}(x_k - \gamma \nabla F(x_k))$$

$$\text{prox}_{\gamma R}(x) = \underset{y}{\operatorname{argmin}} R(y) + \frac{1}{2\gamma} \|y - x\|^2$$

$\text{prox}_{\gamma R}(x)$  easy to compute in some important cases

e.g. explicit expression for  $R = \|\cdot\|_1$  (soft-thresholding)

**Identification:** Beyond convergence

after a finite moment  $K$ , all iterates  $x_k$  ( $k \geq K$ ) lie in an active set  $M$

- Used in e.g. safe screening [El Gahoui '12] [Salmon et al'19] [Sun et al'20]
- We even have bounds on  $K$  [Sun et al '19]
- When the problem is well-posed e.g. [Wright '96], [Lewis Drusvyatskiy '13]

## Enlarged activity identification

Theorem (Fadili, M., Peyré '18)

*Under convergence assumptions, if  $R$  is mirror-stratifiable, then for  $k \geq K$*

$$M_{x^*} \leq M_{x_k} \leq \mathcal{J}_{R^*}(M_{-\nabla F(x^*)}^*)$$

- Optimality condition  $-\nabla F(x^*) \in \partial R(x^*)$

In the **non**-degenerate case:  $-\nabla F(x^*) \in \text{ri}(\partial R(x^*))$

we have exact identification  $M_{x^*} = M_{x_k}$  ( $= \mathcal{J}_{R^*}(M_{-\nabla F(x^*)}^*)$ ) [Liang et al 15]

## Enlarged activity identification

Theorem (Fadili, M., Peyré '18)

*Under convergence assumptions, if  $R$  is mirror-stratifiable, then for  $k \geq K$*

$$M_{x^*} \leq M_{x_k} \leq \mathcal{J}_{R^*}(M_{-\nabla F(x^*)}^*)$$

- Optimality condition  $-\nabla F(x^*) \in \partial R(x^*)$

In the **non**-degenerate case:  $-\nabla F(x^*) \in \text{ri}(\partial R(x^*))$

we have exact identification  $M_{x^*} = M_{x_k}$  ( $= \mathcal{J}_{R^*}(M_{-\nabla F(x^*)}^*)$ ) [Liang et al 15]

- In the general case:  $\delta$  quantifies the degeneracy of the problem

$$\delta = \dim(\mathcal{J}_{R^*}(M_{-\nabla F(x^*)}^*)) - \dim(M_{x^*})$$

$\delta = 0$  : well-posedness (fast convergence and identification)

$\delta$  large : strong degeneracy (slow convergence and identification)

- Note:  $\delta$  and  $K$  are not computable beforehand in general...

## Illustration with nuclear norm

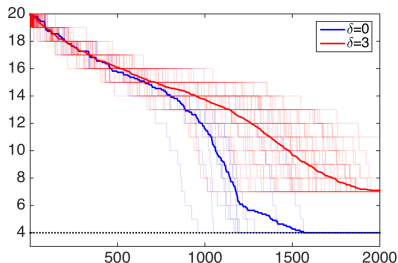
Matrix least-squares regularized by nuclear norm  $\|X\|_* = \|\sigma(X)\|_1$

$$\min_{X \in \mathbb{R}^{d=m \times m}} \frac{1}{2} \|A(X) - y\|^2 + \lambda \|X\|_*$$

Generate many random problems (with  $m = 20$  and  $n = 300$ ), solve them

Select those with  $\text{rank}(X^*) = 4$  and  $\delta = 0$  or  $3$  ( $\delta = \#\{i : |\sigma_i(U^*)| = 1\} - \text{rank}(X^*)$ )

Plot the decrease of  $\text{rank}(X_k)$  with  $X_{k+1} = \text{prox}_{\gamma \|\cdot\|_*}(X_k - \gamma A^*(A(X_k) - y))$



$\delta = 0$ : well-posed vs.  $\delta = 3$ : degenerate

# Outline

- 1 General stability of regularized problems
- 2 Enlarged identification of proximal algorithms
- 3 Application: communication-efficient federated learning
- 4 Application: model consistency for regularized least-squares

## Basic distributed learning set-up

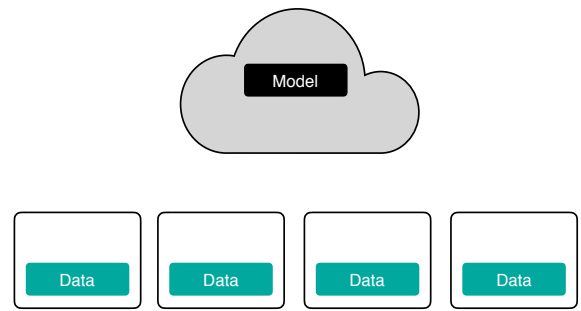
(Standard) centralized learning



- Data  $(a_j, y_j)_{j=1, \dots, n}$ , prediction function  $h(\cdot, x)$ , model parameters  $x \in \mathbb{R}^d$

## Basic distributed learning set-up

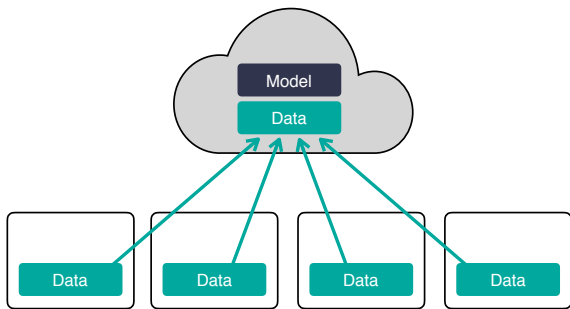
(Standard) centralized learning



- Data  $(a_j, y_j)_{j=1, \dots, n}$ , prediction function  $h(\cdot, x)$ , model parameters  $x \in \mathbb{R}^d$

## Basic distributed learning set-up

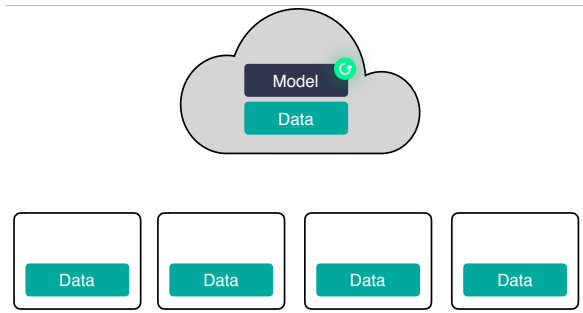
(Standard) centralized learning



- Data  $(a_j, y_j)_{j=1, \dots, n}$ , prediction function  $h(\cdot, x)$ , model parameters  $x \in \mathbb{R}^d$

## Basic distributed learning set-up

(Standard) centralized learning



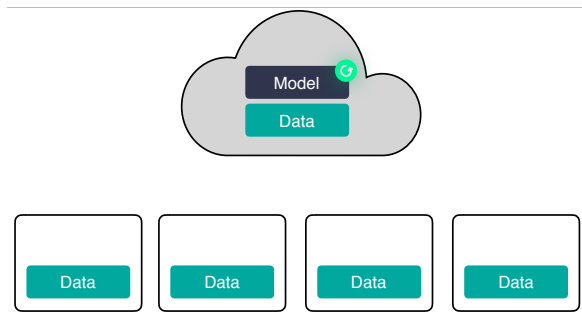
- Data  $(a_j, y_j)_{j=1, \dots, n}$ , prediction function  $h(\cdot, x)$ , model parameters  $x \in \mathbb{R}^d$
- Empirical risk minimization

$$\min_{x \in \mathbb{R}^d} \quad \frac{1}{n} \sum_{j=1}^n \ell(y_j, h(a_j, x)) \quad + \quad \lambda R(x)$$

# Basic distributed learning set-up

(Standard) centralized learning

- needs of lot of storage ☹️
- is highly privacy invasive ☹️

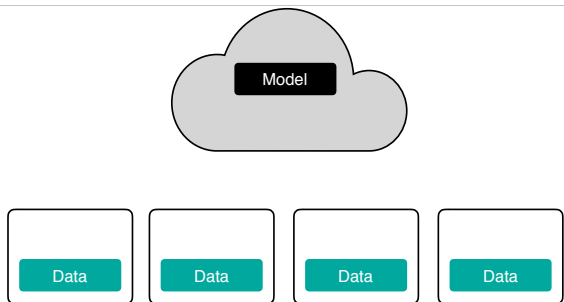


- Data  $(a_j, y_j)_{j=1, \dots, n}$ , prediction function  $h(\cdot, x)$ , model parameters  $x \in \mathbb{R}^d$
- Empirical risk minimization

$$\min_{x \in \mathbb{R}^d} \quad \frac{1}{n} \sum_{j=1}^n \ell(y_j, h(a_j, x)) \quad + \quad \lambda R(x)$$

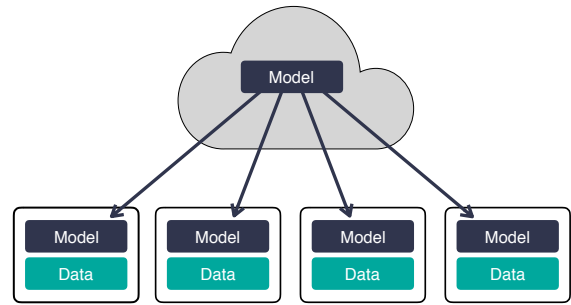
## Move the model, not the data !

Collaborative/Federative learning (introduction of Aurélien's talk this morning)



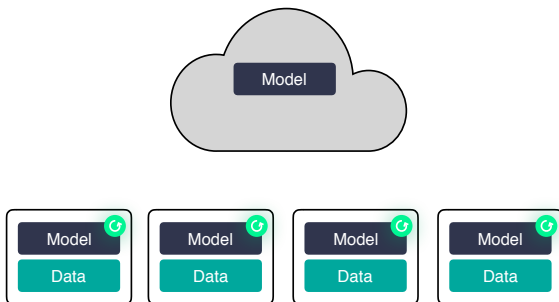
Move the model, not the data !

Collaborative/Federative learning (introduction of Aurélien's talk this morning)



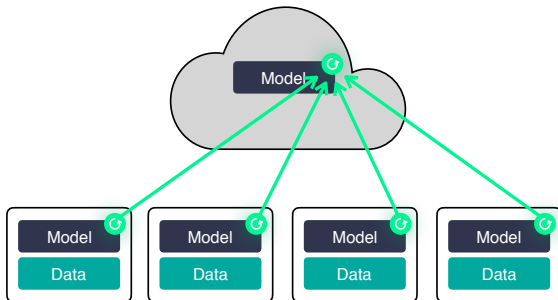
Move the model, not the data !

Collaborative/Federative learning (introduction of Aurélien's talk this morning)



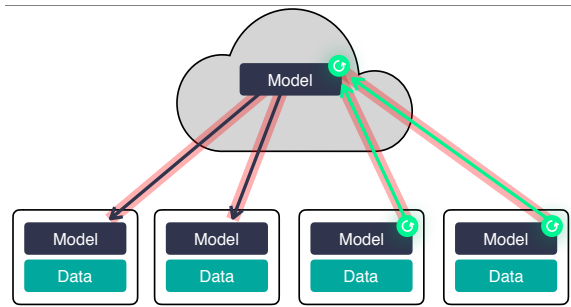
Move the model, not the data !

Collaborative/Federative learning (introduction of Aurélien's talk this morning)



## Move the model, not the data !

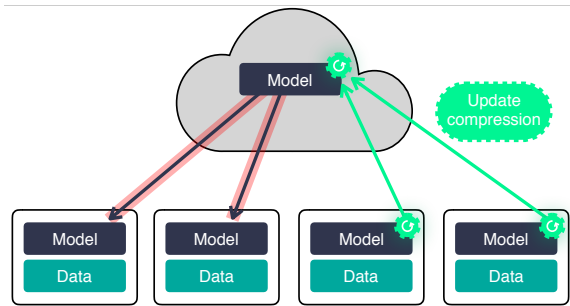
### Collaborative/Federative learning (introduction of Aurélien's talk this morning)



- Communication is the bottleneck 😞
- We need compression ! Mikael talk, yesterday morning (?)

## Move the model, not the data !

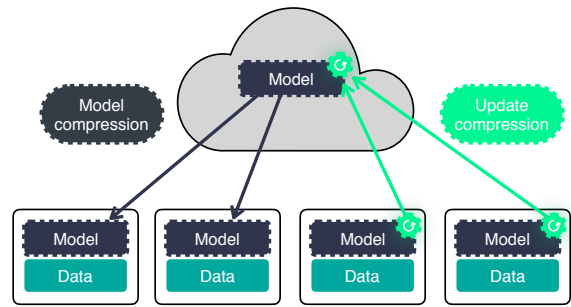
### Collaborative/Federative learning (introduction of Aurélien's talk this morning)



- Communication is the bottleneck 😞
- We need compression ! Mikael talk, yesterday morning (?)

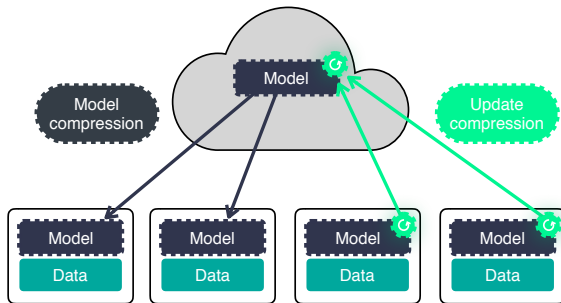
# Move the model, not the data !

## Collaborative/Federative learning (introduction of Aurélien's talk this morning)



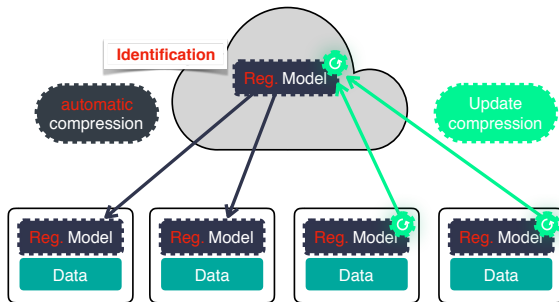
- Communication is the bottleneck 😞
- We need compression ! Mikael talk, yesterday morning (?)
- Many compression techniques... recall Martin's talk yesterday afternoon
- Let's discuss another one, complementary to existing ones

## Application of identification to federated learning



## Application of identification to federated learning

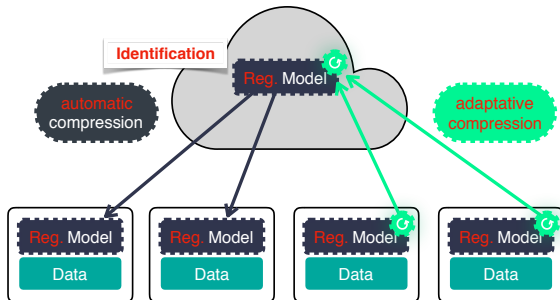
With nonsmooth regularizers, identification comes into play



- Observation: identification gives automatic model compression  
e.g. for  $R = \|\cdot\|_1$ , model becomes sparse... just communicate nonzero entries!

## Application of identification to federated learning

With nonsmooth regularizers, identification comes into play



- Observation: identification gives automatic model compression  
e.g. for  $R = \|\cdot\|_1$ , model becomes sparse... just communicate nonzero entries!
- [Grishchenko, Iutzeler, M. '19] uses again identification for update comp.

Project update onto  $M_{x_k} +$  randomly selected  $M$

e.g. for  $R = \|\cdot\|_1$ , select current support + random entries

- Algo with intricate convergence analysis due to non-uniform selection...

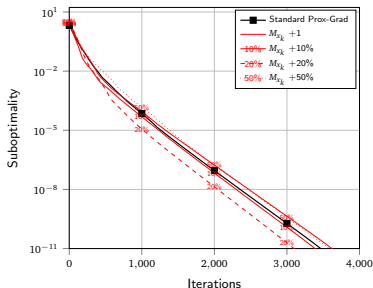
## Illustration of communication-efficient proximal method

On an instance of TV-regularized logistic regression (a1a dataset on 10 machines)

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{j=1}^n \log(1 + \exp(-y_j \langle a_j, x \rangle)) + \lambda \text{TV}(x)$$

$$\text{TV}(x) = \sum_{i=1}^{n-1} |x_{i+1} - x_i|$$

- Comparison of
- Usual distributed proximal-gradient (black)
  - Adaptive distributed proximal-subspace descent (red)
- for different selections  $M_{x_k} + \text{random others}$



# Illustration of communication-efficient proximal method

On an instance of TV-regularized logistic regression (a1a dataset on 10 machines)

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{j=1}^n \log(1 + \exp(-y_j \langle a_j, x \rangle)) + \lambda \text{TV}(x)$$

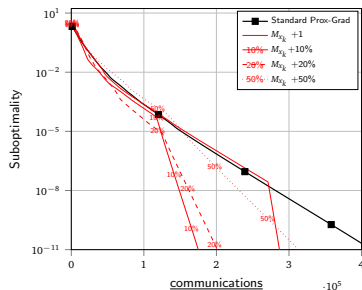
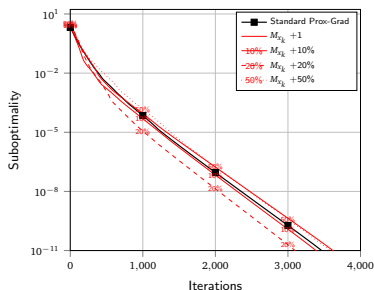
$$\text{TV}(x) = \sum_{i=1}^{n-1} |x_{i+1} - x_i|$$

Total Variation

Comparison of

- Usual distributed proximal-gradient (black)
- Adaptive distributed proximal-subspace descent (red)

for different selections  $M_{x_k} + \text{random others}$



Acceleration... with respect to size of communication

# Illustration of communication-efficient proximal method

On an instance of TV-regularized logistic regression (a1a dataset on 10 machines)

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{j=1}^n \log(1 + \exp(-y_j \langle a_j, x \rangle)) + \lambda \text{TV}(x)$$

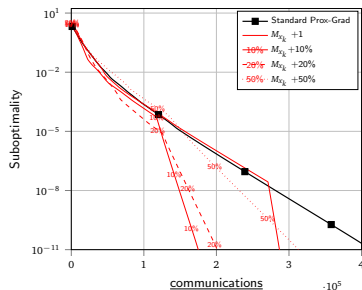
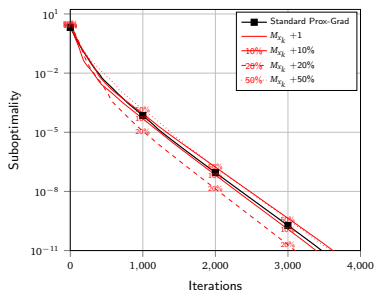
$$\text{TV}(x) = \sum_{i=1}^{n-1} |x_{i+1} - x_i|$$

Total Variation

Comparison of

- Usual distributed proximal-gradient (black)
- Adaptive distributed proximal-subspace descent (red)

for different selections  $M_{x_k} + \text{random others}$



Acceleration... with respect to size of communication

Tradeoff between compression (less comm.) and identification (faster cv)

# Outline

- 1 General stability of regularized problems
- 2 Enlarged identification of proximal algorithms
- 3 Application: communication-efficient federated learning
- 4 Application: model consistency for regularized least-squares

## Supervised learning: model consistency ?

- Assume data  $(a_i, y_i)_{i=1, \dots, n}$  are sampled from linear model

$$y = \langle a, \bar{x} \rangle + \nu \quad \text{with random } (a, \nu)$$

- Structural assumption:  $\bar{x}$  has a low-complexity for  $R$

$$\bar{x} = \operatorname{argmin}_{x \in \mathbb{R}^d} \left\{ R(x) : x \in \operatorname{argmin}_{z \in \mathbb{R}^d} \mathbb{E} \left[ (\langle a, z \rangle - y)^2 \right] \right\}$$

- Regularized least-squares (if  $R = \|\cdot\|_1$ , this is LASSO)

$$\min_{x \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (\langle a_i, x \rangle - y_i)^2 + \lambda_n R(x)$$

- Stochastic (proximal-)gradient algorithms (at iteration  $k$ , pick randomly  $i(k)$ )

$$x_{k+1} = \operatorname{prox}_{\gamma_k \lambda_n R} (x_k - \gamma_k ((\langle a_{i(k)}, x_k \rangle - y_{i(k)}) a_{i(k)} + \varepsilon_k))$$

E.g. SGD, SAGA [Delfazio et al '14], SVRG [Xiao-Zhang '14]

- Do we have model recovery/consistency i.e.  $x_k \in M_{\bar{x}} ?$

(if we have enough observations, i.e. when  $n \rightarrow +\infty$ )

## Enlarged identification of stochastic algorithms

Theorem (Garrigos, Fadili, M., Peyré '19)

Take  $\lambda_n \rightarrow 0$  with  $\lambda_n \sqrt{n/(\log \log n)} \rightarrow +\infty$ . If  $n$  large enough and for

$$\mathbf{x}_{k+1} = \text{prox}_{\gamma_k \lambda_n R} \left( \mathbf{x}_k - \gamma_k \left( (\langle \mathbf{a}_{i(k)}, \mathbf{x}_k \rangle - y_{i(k)}) \mathbf{a}_{i(k)} + \varepsilon_k \right) \right)$$

with mild assumptions on errors  $\varepsilon_k$  and stepsizes  $\gamma_k$ . Then, for  $k$  large, a.s.

$$M_{\bar{\mathbf{x}}} \leq M_{\mathbf{x}_k} \leq \mathcal{J}_{R^*}(M_{\bar{\boldsymbol{\eta}}}^*)$$

with  $\bar{\boldsymbol{\eta}} = \underset{\boldsymbol{\eta} \in \mathbb{R}^p}{\text{argmin}} \left\{ \boldsymbol{\eta}^\top C^\dagger \boldsymbol{\eta} : \boldsymbol{\eta} \in \partial R(\bar{\mathbf{x}}) \cap \text{Im } C \right\}$  and  $C = \mathbb{E}[\mathbf{a}\mathbf{a}^\top]$

Comments:

- key dual object  $\bar{\boldsymbol{\eta}} \in \partial R(\bar{\mathbf{x}})$  [Vaïter et al '16]
- $\lambda_n$  decreases to 0, but not too fast
- SAGA and SVRG satisfy the “mild” assumption [Poon et al '18]
- (Prox-)SGD does not – and does not identify (e.g. [Lee Wright '12])

## Enlarged identification of stochastic algorithms

Theorem (Garrigos, Fadili, M., Peyré '19)

Take  $\lambda_n \rightarrow 0$  with  $\lambda_n \sqrt{n/(\log \log n)} \rightarrow +\infty$ . If  $n$  large enough and for

$$\mathbf{x}_{k+1} = \text{prox}_{\gamma_k \lambda_n R} \left( \mathbf{x}_k - \gamma_k \left( \langle \mathbf{a}_{i(k)}, \mathbf{x}_k \rangle - y_{i(k)} \right) \mathbf{a}_{i(k)} + \varepsilon_k \right)$$

with mild assumptions on errors  $\varepsilon_k$  and stepsizes  $\gamma_k$ . Then, for  $k$  large, a.s.

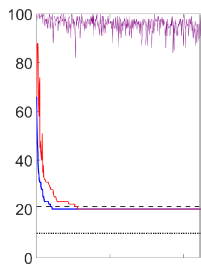
$$M_{\bar{\mathbf{x}}} \leq M_{\mathbf{x}_k} \leq \mathcal{J}_{R^*}(M_{\bar{\boldsymbol{\eta}}}^*)$$

with  $\bar{\boldsymbol{\eta}} = \underset{\boldsymbol{\eta} \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \boldsymbol{\eta}^\top C^\dagger \boldsymbol{\eta} : \boldsymbol{\eta} \in \partial R(\bar{\mathbf{x}}) \cap \operatorname{Im} C \right\}$  and  $C = \mathbb{E}[\mathbf{a}\mathbf{a}^\top]$

Comments:

- key dual object  $\bar{\boldsymbol{\eta}} \in \partial R(\bar{\mathbf{x}})$  [Vaiter et al '16]
- $\lambda_n$  decreases to 0, but not too fast
- SAGA and SVRG satisfy the “mild” assumption [Poon et al '18]
- (Prox-)SGD does not – and does not identify (e.g. [Lee Wright '12])

(on a LASSO instance)



## Conclusion

Take-home message: identification often holds... and can be used

- Enlarged identification results (explaining observed phenomena)
- Better understanding of optim. algos (beyond convergence)
- Sparsify communications by adaptative dimension reduction

# Conclusion

Take-home message: identification often holds... and can be used

- Enlarged identification results (explaining observed phenomena)
- Better understanding of optim. algos (beyond convergence)
- Sparsify communications by adaptative dimension reduction

Extensions, on-going work

- Many possible refinements of sensitivity results  
other data fidelity terms, a priori control on strata dimension, explaining transition curves...
- Use identification to accelerate convergence  
interplay between identification and acceleration (PhD of Gilles Bareilles)
- Subspace descent algorithms generalizing coordinate descent  
"coordinate" descent for nonseparable functions → Franck's talk tomorrow

## Conclusion

Take-home message: identification often holds... and can be used

- Enlarged identification results (explaining observed phenomena)
- Better understanding of optim. algos (beyond convergence)
- Sparsify communications by adaptative dimension reduction

## Extensions, on-going work

- Many possible refinements of sensitivity results  
other data fidelity terms, a priori control on strata dimension, explaining transition curves...
- Use identification to accelerate convergence  
interplay between identification and acceleration (PhD of Gilles Bareilles)
- Subspace descent algorithms generalizing coordinate descent  
"coordinate" descent for nonseparable functions → Franck's talk tomorrow

thanks !!