

An Overview of Stochastic Methods for Solving Optimization Problems

Émilie Chouzenoux

Laboratoire d'Informatique Gaspard Monge - CNRS
Univ. Paris-Est Marne-la-Vallée, France

26 Nov. 15



Introduction

STOCHASTIC PROBLEM

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \mathbb{E}(\varphi_j(\mathbf{h}_j^\top \mathbf{x}, y_j)) + g(\mathbf{D}\mathbf{x})$$

where $j \in \mathbb{N}^*$, $\mathbf{h}_j \in \mathbb{R}^N$, $y_j \in \mathbb{R}$, $\varphi_j: \mathbb{R} \times \mathbb{R} \rightarrow]-\infty, +\infty]$ is a loss function, and $g \circ \mathbf{D}$ is a regularization function, with $g: \mathbb{R}^P \rightarrow]-\infty, +\infty]$ and $\mathbf{D} \in \mathbb{R}^{P \times N}$.

Introduction

STOCHASTIC PROBLEM

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \mathbb{E}(\varphi_j(\mathbf{h}_j^\top \mathbf{x}, y_j)) + g(\mathbf{D}\mathbf{x})$$

where $j \in \mathbb{N}^*$, $\mathbf{h}_j \in \mathbb{R}^N$, $y_j \in \mathbb{R}$, $\varphi_j: \mathbb{R} \times \mathbb{R} \rightarrow]-\infty, +\infty]$ is a loss function, and $g \circ \mathbf{D}$ is a regularization function, with $g: \mathbb{R}^P \rightarrow]-\infty, +\infty]$ and $\mathbf{D} \in \mathbb{R}^{P \times N}$.

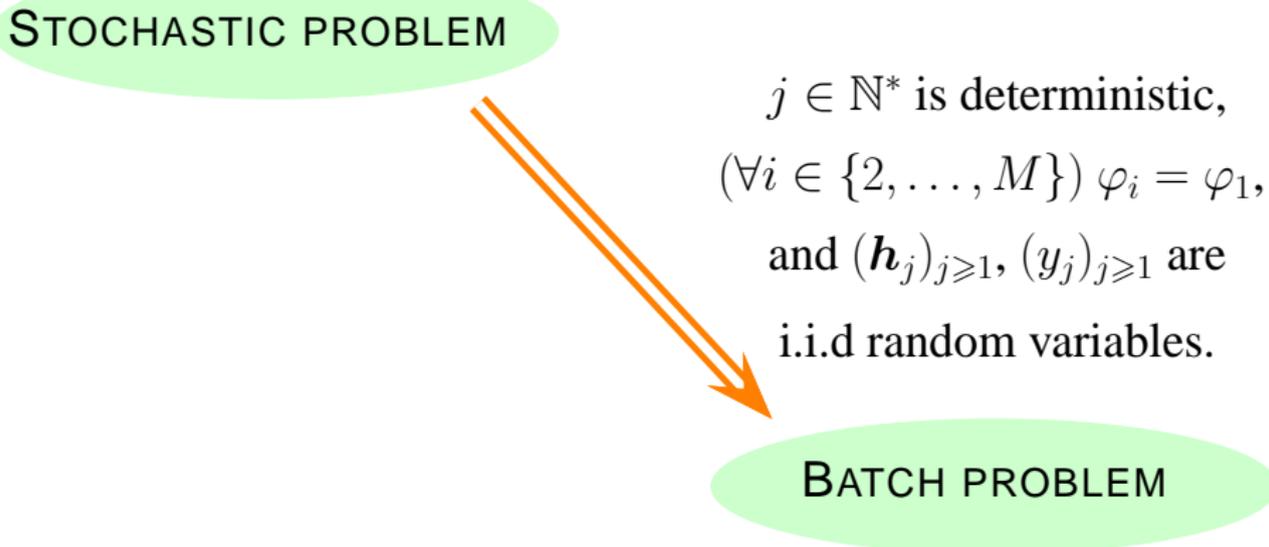
BATCH PROBLEM

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \frac{1}{M} \sum_{i=1}^M \varphi_i(\mathbf{h}_i^\top \mathbf{x}, y_i) + g(\mathbf{D}\mathbf{x})$$

where, for all $i \in \{1, \dots, M\}$, $\varphi_i: \mathbb{R} \times \mathbb{R} \rightarrow]-\infty, +\infty]$, $\mathbf{h}_i \in \mathbb{R}^N$ and $y_i \in \mathbb{R}$.

Link between stochastic and batch problems

STOCHASTIC PROBLEM



$j \in \mathbb{N}^*$ is deterministic,
 $(\forall i \in \{2, \dots, M\}) \varphi_i = \varphi_1$,
and $(\mathbf{h}_j)_{j \geq 1}, (y_j)_{j \geq 1}$ are
i.i.d random variables.

BATCH PROBLEM

Link between stochastic and batch problems

STOCHASTIC PROBLEM

\mathbf{y} and \mathbf{H} are deterministic,
and j is uniformly distributed
over $\{1, \dots, M\}$.

BATCH PROBLEM



Introduction

NUMEROUS EXAMPLES:

- ▶ supervised classification
- ▶ inverse problems
- ▶ system identification, channel equalization
- ▶ linear prediction/interpolation
- ▶ echo cancellation, interference removal
- ▶ ...

In the context of **large scale problems**, how to find an optimization algorithm able to deliver a reliable numerical solution in a **reasonable time**, with **low memory requirement**?

Outline

- * FUNDAMENTAL TOOLS IN CONVEX ANALYSIS
- * OPTIMIZATION ALGORITHMS FOR SOLVING STOCHASTIC PROBLEM
 - ▶ Stochastic forward-backward algorithm
 - ▶ A brief focus on sparse adaptive filtering
- * STOCHASTIC ALGORITHMS FOR SOLVING BATCH PROBLEM
 - ▶ Incremental gradient algorithms
 - ▶ Block coordinate approaches

Fundamental tools in convex analysis

Notation and definitions

Let $f: \mathbb{R}^N \rightarrow]-\infty, +\infty]$.

- ▶ The **domain** of function f is

$$\text{dom } f = \{ \mathbf{x} \in \mathbb{R}^N \mid f(\mathbf{x}) < +\infty \}$$

If $\text{dom } f \neq \emptyset$, function f is said to be **proper**.

- ▶ Function f is **convex** if

$$(\forall (\mathbf{x}, \mathbf{y}) \in (\mathbb{R}^N)^2)(\forall \lambda \in [0, 1])$$

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}).$$

- ▶ Function f is **lower semi-continuous** (lsc) on \mathbb{R}^N if, for all $\mathbf{x} \in \mathbb{R}^N$, for all sequence $(\mathbf{x}_k)_{k \in \mathbb{N}}$ of \mathbb{R}^N ,

$$\mathbf{x}_k \longrightarrow \mathbf{x} \quad \Rightarrow \quad \liminf f(\mathbf{x}_k) \geq f(\mathbf{x}).$$

Notation and definitions

Let $f: \mathbb{R}^N \rightarrow]-\infty, +\infty]$. Function f is said ν -strongly convex if

$$(\forall (\mathbf{x}, \mathbf{y}) \in (\mathbb{R}^N)^2)(\forall \lambda \in [0, 1])$$

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) - \frac{1}{2} \nu \lambda (1 - \lambda) \|\mathbf{x} - \mathbf{y}\|^2,$$

with $\nu \in]0, +\infty[$.

Notation and definitions

Let $f: \mathbb{R}^N \rightarrow]-\infty, +\infty]$. Function f is said **ν -strongly convex** if

$$(\forall (\mathbf{x}, \mathbf{y}) \in (\mathbb{R}^N)^2)(\forall \lambda \in [0, 1])$$

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) - \frac{1}{2} \nu \lambda (1 - \lambda) \|\mathbf{x} - \mathbf{y}\|^2,$$

with $\nu \in]0, +\infty[$.

Let $f: \mathbb{R}^N \rightarrow]-\infty, +\infty[$. Function f is said **β -Lipschitz differentiable** if it is differentiable over \mathbb{R}^N and its gradient fulfills

$$(\forall (\mathbf{x}, \mathbf{y}) \in (\mathbb{R}^N)^2) \quad \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq \beta \|\mathbf{x} - \mathbf{y}\|,$$

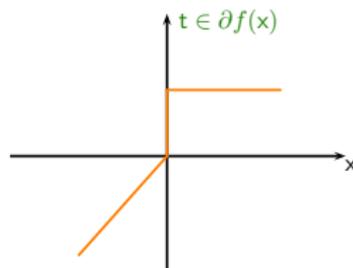
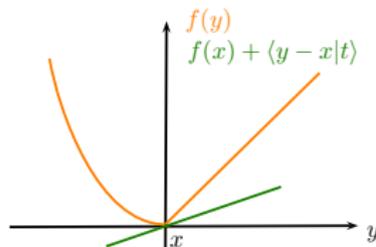
with $\beta \in]0, +\infty[$.

Subdifferential

The **subdifferential** of a convex function $f: \mathbb{R}^N \rightarrow]-\infty, +\infty]$ at x is the set

$$\partial f(x) = \{t \in \mathbb{R}^N \mid (\forall y \in \mathbb{R}^N) f(y) \geq f(x) + \langle t \mid y - x \rangle\}$$

An element t of $\partial f(x)$ is called **a subgradient** of f at x .



- If f is differentiable at $x \in \mathbb{R}^N$ then $\partial f(x) = \{\nabla f(x)\}$.

Proximity operator

Let $f : \mathbb{R}^N \mapsto] - \infty, +\infty]$ a proper, convex, l.s.c function.

CHARACTERIZATION OF PROXIMITY OPERATOR

$$(\forall \mathbf{x} \in \mathbb{R}^N) \quad \hat{\mathbf{y}} = \text{prox}_f(\mathbf{x}) \Leftrightarrow \mathbf{x} - \hat{\mathbf{y}} \in \partial f(\hat{\mathbf{y}}).$$

The proximity operator $\text{prox}_f(\mathbf{x})$ of f at $\mathbf{x} \in \mathbb{R}^N$ is the unique vector $\hat{\mathbf{y}} \in \mathbb{R}^N$ such that

$$f(\hat{\mathbf{y}}) + \frac{1}{2} \|\hat{\mathbf{y}} - \mathbf{x}\|^2 = \inf_{\mathbf{y} \in \mathbb{R}^N} f(\mathbf{y}) + \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

Properties of proximal operator

	$f(\mathbf{x})$	$\text{prox}_f(\mathbf{x})$
translation $\mathbf{z} \in \mathbb{R}^N$	$f(\mathbf{x} - \mathbf{z})$	$\mathbf{z} + \text{prox}_f(\mathbf{x} - \mathbf{z})$
quadratic perturbation $\mathbf{z} \in \mathbb{R}^N, \alpha > 0, \gamma \in \mathbb{R}$	$f(\mathbf{x}) + \alpha \ \mathbf{x}\ ^2/2 + \langle \mathbf{x} \mathbf{z} \rangle + \gamma$	$\text{prox}_{\frac{f}{\alpha+1}}\left(\frac{\mathbf{x}-\mathbf{z}}{\alpha+1}\right)$
scaling $\rho \in \mathbb{R}^*$	$f(\rho\mathbf{x})$	$\frac{1}{\rho} \text{prox}_{\rho^2 f}(\rho\mathbf{x})$
quadratic function $\mathbf{L} \in \mathbb{R}^{M \times N}, \gamma > 0, \mathbf{z} \in \mathbb{R}^M$	$\gamma \ \mathbf{L}\mathbf{x} - \mathbf{z}\ ^2/2$	$(\text{Id} + \gamma \mathbf{L}\mathbf{L}^*)^{-1}(\mathbf{x} - \gamma \mathbf{L}^* \mathbf{z})$
semi-unitary transform $\mathbf{L} \in \mathbb{R}^{M \times N}, \mathbf{L}\mathbf{L}^* = \mu \text{Id}, \mu > 0$	$f(\mathbf{L}\mathbf{x})$	$\mathbf{x} - \mu^{-1} \mathbf{L}^*(\mathbf{x} - \text{prox}_{\mu f}(\mathbf{L}\mathbf{x}))$
reflexion	$f(-\mathbf{x})$	$-\text{prox}_f(-\mathbf{x})$
separability	$\sum_{i=1}^N \varphi_i(x^{(i)})$ $\mathbf{x} = (x^{(i)})_{1 \leq i \leq N}$	$\left(\text{prox}_{\varphi_i}(x^{(i)})\right)_{1 \leq i \leq N}$
indicator function	$\iota_C(\mathbf{x})$	$\text{P}_C(\mathbf{x})$
support function	$\iota_C^*(\mathbf{x}) = \sigma_C(\mathbf{x})$	$\mathbf{x} - \text{P}_C(\mathbf{x})$

Optimization algorithms for solving stochastic problem

Stochastic forward-backward algorithm

STOCHASTIC PROBLEM

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \mathbb{E}(\varphi_j(\mathbf{h}_j^\top \mathbf{x}, y_j)) + g(\mathbf{D}\mathbf{x})$$

⇒ At each iteration $j \geq 1$, assume that an estimate \mathbf{u}_j of the gradient of $\Phi(\cdot) = \mathbb{E}(\varphi_j(\mathbf{h}_j^\top \cdot, y_j))$ at \mathbf{x}_j is available.

Stochastic forward-backward algorithm

STOCHASTIC PROBLEM

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \mathbb{E}(\varphi_j(\mathbf{h}_j^\top \mathbf{x}, y_j)) + g(D\mathbf{x})$$

⇒ At each iteration $j \geq 1$, assume that an estimate \mathbf{u}_j of the **gradient** of $\Phi(\cdot) = \mathbb{E}(\varphi_j(\mathbf{h}_j^\top \cdot, y_j))$ at \mathbf{x}_j is available.

The SFB algorithm reads:

$$\begin{aligned} & (\gamma_j)_{j \geq 1} \in]0, +\infty[, (\lambda_j)_{j \geq 1} \in]0, 1] \\ & \text{for } j = 1, 2, \dots \\ & \left[\begin{array}{l} \mathbf{z}_j = \text{prox}_{\gamma_j g \circ D}(\mathbf{x}_j - \gamma_j \mathbf{u}_j) \\ \mathbf{x}_{j+1} = (1 - \lambda_j) \mathbf{x}_j + \lambda_j \mathbf{z}_j \end{array} \right. \end{aligned}$$

- ▶ When $g \equiv 0$, the stochastic gradient descent (SGD) algorithm is recovered.

Convergence theorem [Rosasco *et al.*, 2014]

Let $F \neq \emptyset$ denote the set of minimizers of $\Phi + g \circ D$. Assume that:

- (i) Φ has a β -Lipschitzian gradient with $\beta \in]0, +\infty[$, g is a proper, lower-semicontinuous convex function, and $\Phi + g \circ D$ is strongly convex.
- (ii) For every $j \geq 1$,

$$\mathbb{E}(\{\|\mathbf{u}_j\|^2\}) < +\infty, \quad \mathbb{E}\{\mathbf{u}_j \mid \mathcal{X}_{j-1}\} = \nabla\Phi(\mathbf{x}_j),$$

$$\mathbb{E}\{\|\mathbf{u}_j - \nabla\Phi(\mathbf{x}_j)\|^2 \mid \mathcal{X}_{j-1}\} \leq \sigma^2(1 + \alpha_j \|\nabla\Phi(\mathbf{x}_j)\|^2)$$

where $\mathcal{X}_j = (y_i, \mathbf{h}_i)_{1 \leq i \leq j}$, and α_j and σ are positive values such that $\gamma_j \leq (2 - \epsilon)/(\beta(1 + 2\sigma^2\alpha_j))$ with $\epsilon > 0$.

- (iii) We have

$$\sum_{j \geq 1} \lambda_j \gamma_j = +\infty \quad \text{and} \quad \sum_{j \geq 1} \chi_j^2 < +\infty$$

where, for every $j \geq 1$, $\chi_j^2 = \lambda_j \gamma_j^2 (1 + 2\alpha_j \|\nabla\Phi(\bar{\mathbf{x}})\|^2)$ and $\bar{\mathbf{x}} \in F$.

Then, $(\mathbf{x}_j)_{j \geq 1}$ converges almost surely to an element of F .

Bibliographical remarks

RELATED APPROACHES

- ▶ Methods relying on subgradient steps [Shalev-Shwartz *et al.*, 2007],
- ▶ Regularized dual averaging methods [Xiao, 2010],
- ▶ Composite mirror descent methods [Duchi *et al.*, 2010].

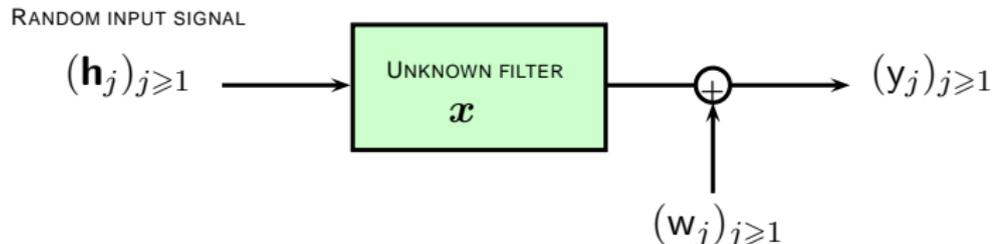
WHAT IF PROX OF $g \circ D$ IS NOT SIMPLE?

- ▶ Stochastic proximal averaging strategy [Zhong *et al.*, 2014],
- ▶ Conditional gradient (\sim Frank-Wolfe) techniques [Lafond, 2015],
- ▶ Stochastic ADMM [Ouyang *et al.*, 2013],
- ▶ Block alternating strategy [Xu *et al.*, 2014],
- ▶ Stochastic proximal primal-dual methods (also for varying g) [Combettes *et al.*, 2015].

HOW TO ACCELERATE CONVERGENCE?

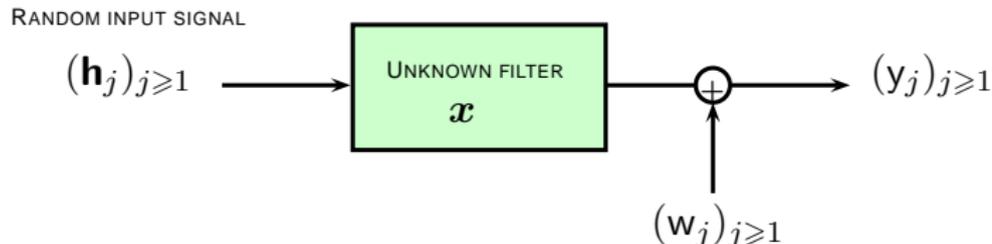
- ▶ Subspace acceleration techniques [Hu *et al.*, 2009][Atchadé *et al.*, 2014],
- ▶ Preconditioning techniques [Duchi *et al.*, 2011],
- ▶ Mixing both strategies (smooth case) [Chouzenoux *et al.*, 2014].

A brief focus on sparse adaptive filtering



⇒ Previous stochastic problem, with $(\forall j \geq 1) \varphi_j(\mathbf{h}_j^\top \mathbf{x}, y_j) = (\mathbf{h}_j^\top \mathbf{x} - y_j)^2$.

A brief focus on sparse adaptive filtering



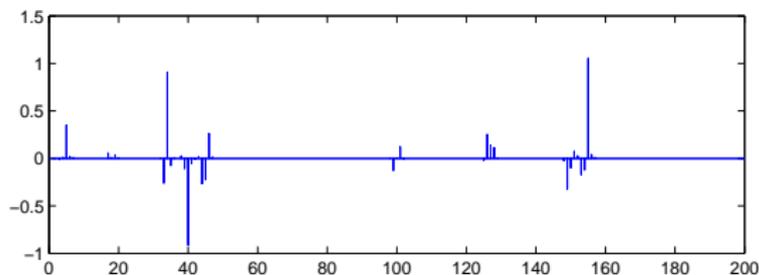
⇒ Previous stochastic problem, with $(\forall j \geq 1) \varphi_j(\mathbf{h}_j^\top \mathbf{x}, y_j) = (\mathbf{h}_j^\top \mathbf{x} - y_j)^2$.

EXISTING WORKS IN CASE OF SPARSE PRIOR:

- * Proportionate least mean square methods (\sim Preconditioned SGD) [Paleologu *et al.*, 2010],
- * Zero-attracting algorithms (\sim subgradient descent) [Chen *et al.*, 2010],
- * Proximal-like algorithms: SFB [Yamagashi *et al.*, 2011] or primal-dual approach [Ono *et al.*, 2013],
- * Penalized versions of recursive least squares [Angelosante *et al.*, 2011],
- * Over-relaxed projection algorithms [Kopsinis *et al.*, 2011],
- * Time-varying filters \rightsquigarrow affine projection strategy (\sim mini-batch in machine learning) [Markus *et al.*, 2014].

Simulation results

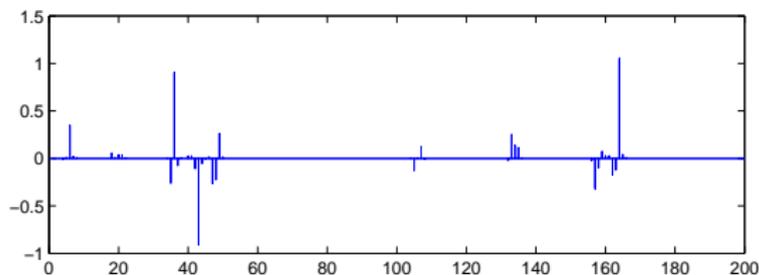
- x : Time-variant linear system with 200 sparse coefficients,
- h : Input sequence of 5000 random independent variables uniformly distributed on $\{-1, +1\}$,
- w : White Gaussian noise with zero mean and variance 0.05.



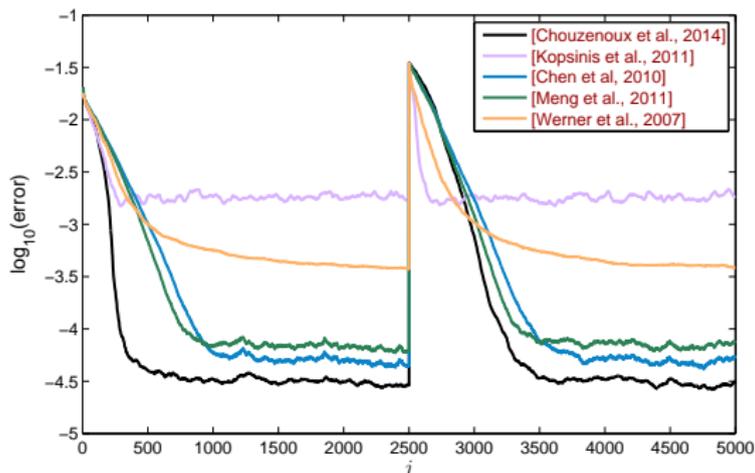
Values of the coefficients of the true sparse filter x for $1 \leq j \leq 2500$

Simulation results

- x : Time-variant linear system with 200 sparse coefficients,
- h : Input sequence of 5000 random independent variables uniformly distributed on $\{-1, +1\}$,
- w : White Gaussian noise with zero mean and variance 0.05.



Values of the coefficients of the true sparse filter x for $2501 \leq j \leq 5000$



Estimation error along time, for various sparse adaptive filtering strategies

- ▶ The parameters of each tested method (forgetting factor, stepsize, regularization weight, affine projection blocksize) are optimized manually,
- ▶ The Stochastic Majorize-Minimize Memory gradient (S3MG) algorithm from [Chouzenoux *et al.*, 2014] leads to a minimal estimation error, while benefiting from good tracking properties.

Stochastic algorithms for solving batch problem

Incremental gradient algorithms

BATCH PROBLEM

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \frac{1}{M} \sum_{i=1}^M \varphi_i(\mathbf{h}_i^\top \mathbf{x}, y_i) + g(\mathbf{D}\mathbf{x})$$

⇒ At each iteration $n \geq 0$, some $j_n \in \{1, \dots, M\}$ is randomly chosen, and only the gradient of $\varphi_{j_n}(\mathbf{h}_{j_n}^\top \cdot, y_{j_n})$ at \mathbf{x}_n is computed.

Incremental gradient algorithms

BATCH PROBLEM

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \frac{1}{M} \sum_{i=1}^M \varphi_i(\mathbf{h}_i^\top \mathbf{x}, y_i) + g(\mathbf{D}\mathbf{x})$$

⇒ At each iteration $n \geq 0$, some $j_n \in \{1, \dots, M\}$ is randomly chosen, and only the gradient of $\varphi_{j_n}(\mathbf{h}_{j_n}^\top \cdot, y_{j_n})$ at \mathbf{x}_n is computed.

For instance, the SAGA algorithm [Defazio *et al.*, 2014] reads:

$\gamma \in]0, +\infty[$, and $(\forall i \in \{1, \dots, M\}) \mathbf{z}_{i,0} = \mathbf{x}_0 \in \mathbb{R}^N$.
for $n = 0, 1, \dots$

$$\left[\begin{array}{l} \text{Select randomly } j_n \in \{1, \dots, M\}, \\ \mathbf{u}_n = \mathbf{h}_{j_n} \nabla \varphi_{j_n}(\mathbf{h}_{j_n}^\top \mathbf{x}_n, y_{j_n}) - \mathbf{h}_{j_n} \nabla \varphi_{j_n}(\mathbf{h}_{j_n}^\top \mathbf{z}_{j_n,n}, y_{j_n}) \\ \quad + \frac{1}{M} \sum_{i=1}^M \mathbf{h}_i \nabla \varphi_i(\mathbf{h}_i^\top \mathbf{z}_{i,n}, y_i) \\ \mathbf{x}_{n+1} = \text{prox}_{\gamma g \circ \mathbf{D}}(\mathbf{x}_n - \gamma \mathbf{u}_n) \\ \mathbf{z}_{j_n,n+1} = \mathbf{x}_{n+1}, \text{ and } (\forall i \in \{1, \dots, M\}) \mathbf{z}_{i,n+1} = \mathbf{z}_{i,n} \end{array} \right.$$

Convergence theorem [Defazio *et al.*, 2014]

Let $\Phi(\cdot) = \frac{1}{M} \sum_{i=1}^M \varphi_i(\mathbf{h}_i^\top \cdot, y_i)$. Denote by $F \neq \emptyset$ the set of minimizers of $\Phi + g \circ \mathbf{D}$. If:

- (i) Φ is convex, β -Lipschitz differentiable on \mathbb{R}^N , and g is proper, lower-semicontinuous convex on \mathbb{R}^N ,
- (ii) For every $n \in \mathbb{N}$, j_n is drawn from an i.i.d. uniform distribution on $\{1, \dots, M\}$,

Then, for $\gamma = 1/3\beta$, for $n \in \mathbb{N}^*$,

$$\mathbb{E}((\Phi + g \circ \mathbf{D})(\bar{\mathbf{x}}_n)) - (\Phi + g \circ \mathbf{D})(\hat{\mathbf{x}}) \leq \frac{4M}{n} \left(\frac{2\beta}{M} \|\mathbf{x}_0 - \hat{\mathbf{x}}\|^2 + \Phi(\mathbf{x}_0) - \nabla\Phi(\hat{\mathbf{x}})^\top (\mathbf{x}_0 - \hat{\mathbf{x}}) - \Phi(\hat{\mathbf{x}}) \right),$$

where $\hat{\mathbf{x}} \in F$ and $\bar{\mathbf{x}}_n = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j$.

If, additionally, Φ is ν -strongly convex then, for $\gamma = 1/(2(\nu M + \beta))$,

$$\mathbb{E}(\|\mathbf{x}_n - \hat{\mathbf{x}}\|^2) \leq \left(1 - \frac{\nu}{\gamma}\right)^n (\|\mathbf{x}_0 - \hat{\mathbf{x}}\|^2 + 2\gamma M (\Phi(\mathbf{x}_0) - \nabla\Phi(\hat{\mathbf{x}})^\top (\mathbf{x}_0 - \hat{\mathbf{x}}) - \Phi(\hat{\mathbf{x}}))).$$

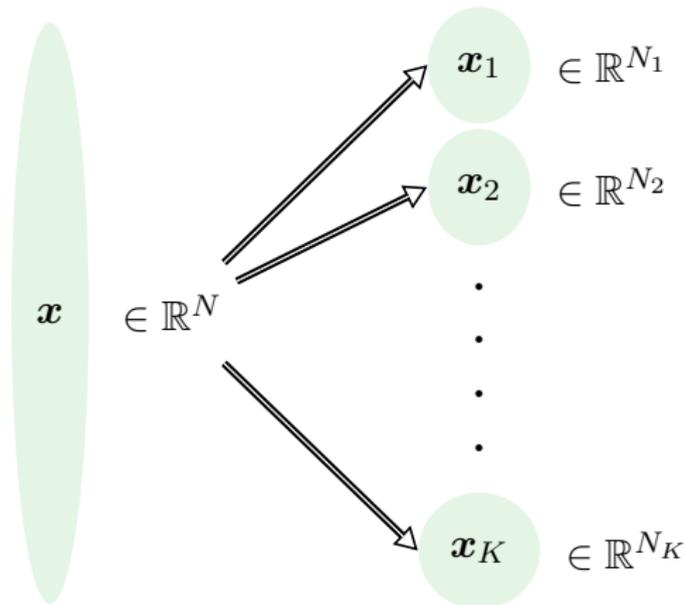
Bibliographical remarks

⇒ Links between stochastic incremental methods existing in the literature:

ALGORITHM	GENERAL IDEA	PROS/CONS	REFS
Standard incremental gradient	$\mathbf{u}_n = \mathbf{h}_{j_n} \nabla \varphi_{j_n}(\mathbf{h}_{j_n}^\top \mathbf{x}_n, y_{j_n})$	simplicity / decreasing stepsize required	[Bertsekas, 2010]
Variance reduction approaches (SVRG / mSGD)	At every $K \geq 0$ iterations, perform a full gradient step (\sim mini-batch strategy)	reduced memory / more gradient evaluations	[Konečný, 2014], [Johnson et al, 2014]
Gradient averaging (SAG / SAGA)	Factor $1/M$ in front of gradient difference term	lower variance / increasing bias (in gradient estimates)	[Schmidt et al, 2014], [Defazio et al, 2014]
Proximal averaging (FINITO)	$\mathbf{x}_{n+1} = \text{prox}_{\gamma g \circ D}(\bar{\mathbf{z}}_n - \gamma \mathbf{u}_n)$ with $\bar{\mathbf{z}}_n$ average of $(\mathbf{z}_{i,n})_{1 \leq i \leq M}$	extra storage cost / less gradient evaluations	[Defazio et al., 2014]
Majorization-Minimization (MISO)	\mathbf{x}_{n+1} minimizer of a majorant function of $\varphi_{j_n}(\mathbf{h}_{j_n}^\top \cdot, y_{j_n}) + g \circ D$ at $\bar{\mathbf{z}}_n$	extra storage cost / less gradient evaluations	[Mairal, 2015]

Block coordinate approaches

► Idea: variable splitting.



Assumption: $g(\mathbf{D}\mathbf{x}) = \sum_{k=1}^K g_{1,k}(\mathbf{x}_k) + g_{2,k}(\mathbf{D}_k\mathbf{x}_k)$ where, for every $k \in \{1, \dots, K\}$, $\mathbf{D}_k \in \mathbb{R}^{P_k \times N_k}$.

Stochastic primal-dual proximal algorithm [Pesquet *et al.*, 2015]

$\tau \in]0, +\infty[, \gamma \in]0, +\infty[$,

for $n = 1, 2, \dots$

for $k = 1, 2, \dots, K$

with probability $\varepsilon_k \in]0, 1]$ do

$$\mathbf{v}_{k,n+1} = (\text{Id} - \text{prox}_{\tau^{-1}g_{2,k}})(\mathbf{v}_{k,n} + \mathbf{D}_k \mathbf{x}_{k,n})$$

$$\mathbf{x}_{k,n+1} = \text{prox}_{\gamma g_{1,k}} \left(\mathbf{x}_{k,n} - \gamma \left(\tau \mathbf{D}_k^\top (2\mathbf{v}_{k+1,n} - \mathbf{v}_{k,n}) + \frac{1}{M} \sum_{i=1}^M \mathbf{h}_{i,k} \nabla \varphi_i \left(\sum_{k'=1}^K \mathbf{h}_{i,k'}^\top \mathbf{x}_{k',n}, y_i \right) \right) \right)$$

otherwise

$$\mathbf{v}_{k,n+1} = \mathbf{v}_{k,n}, \quad \mathbf{x}_{k,n+1} = \mathbf{x}_{k,n}.$$

- ▶ When $g_{2,k} \equiv 0$, the random block coordinate forward-backward algorithm is recovered [Combettes *et al.*, 2015],
- ▶ When $g_{1,k} \equiv 0$ and $g_{2,k} \equiv 0$, the random block coordinate descent algorithm is obtained [Nesterov, 2012].

Convergence theorem [Pesquet *et al.*, 2015]

Set, for every $n \in \mathbb{N}^*$, $\mathcal{X}_n = (\mathbf{x}_{n'}, \mathbf{v}_{n'})_{1 \leq n' \leq n}$.

Let $F \neq \emptyset$ denote the set of minimizers of $\Phi + g \circ D$.

Assume that:

- (i) Φ is convex, β -Lipschitz differentiable on \mathbb{R}^N , g is lower-semicontinuous convex on \mathbb{R}^N ,
- (ii) The blocks activation is performed at each iteration n independently of \mathcal{X}_n , with positive probabilities $(\varepsilon_1, \dots, \varepsilon_K)$,
- (iv) The primal and dual stepsizes (τ, γ) satisfy
$$\frac{1}{\tau} - \gamma \max_{1 \leq k \leq K} \|\mathbf{D}_k\|^2 > \frac{\beta}{2},$$

Then, $(\mathbf{x}_n)_{n \in \mathbb{N}^*}$ converges weakly almost surely to an F -valued random variable.

Bibliographical remarks

CONVERGENCE ANALYSIS

- ▶ Almost sure convergence [Pesquet *et al.*, 2015],
- ▶ Worst case convergence rates [Richtarik *et al.*, 2014] [Necoara *et al.*, 2014] [Lu *et al.*, 2015].

VARIANTS OF THE METHOD

- ▶ Improved convergence conditions in some specific cases [Fercoq *et al.*, 2015],
- ▶ Dual ascent strategies in the strongly convex case (\sim dual forward-backward) [Shalev-Shwartz *et al.*, 2014] [Jaggi *et al.*, 2014] [Qu *et al.*, 2014],
- ▶ Douglas-Rachford/ADMM approaches [Combettes *et al.*, 2015] [Iutzeler *et al.*, 2013],
- ▶ Asynchronous distributed algorithms [Pesquet *et al.*, 2014] [Bianchi *et al.*, 2014].

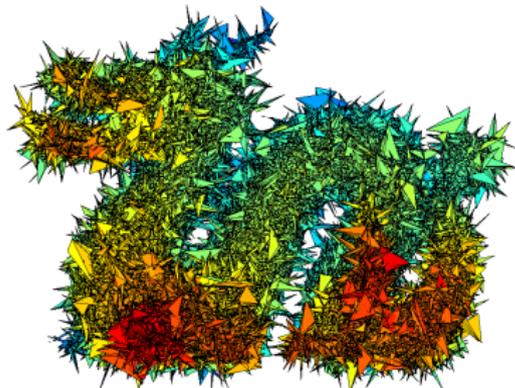
⇒ Dual ascent strategies and asynchronous distributed methods are closely related to **incremental gradient algorithms**.

Simulation results

(ANR GRAPHSIP)



Original mesh, $N = 100250$.

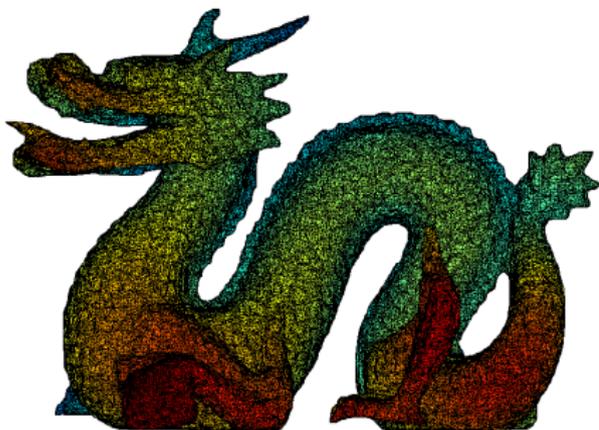


Noisy mesh, $\text{MSE} = 2.89 \times 10^{-6}$.

Goal: Restore the nodes positions of an original mesh corrupted through an additive i.i.d. zero-mean Gaussian mixture noise model,

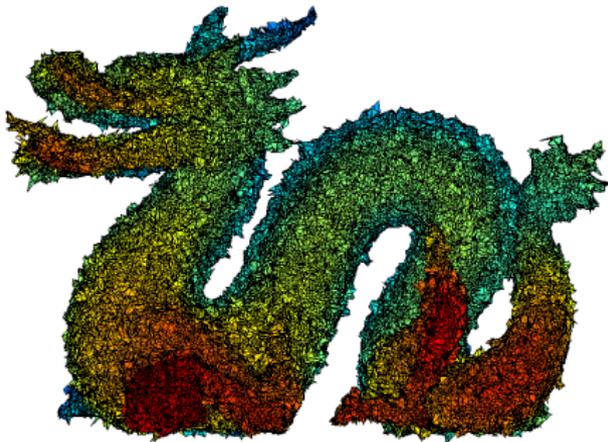
Limited memory available \Rightarrow The mesh is decomposed into K/r non-overlapping blocks with size $r \leq K$, and ϵ is such that only one block is updated at each iteration.

- Reconstruction results using the stochastic primal-dual proximal algorithm for 3D mesh denoising from [Repetti *et al.*, 2015]:



Proposed reconstruction

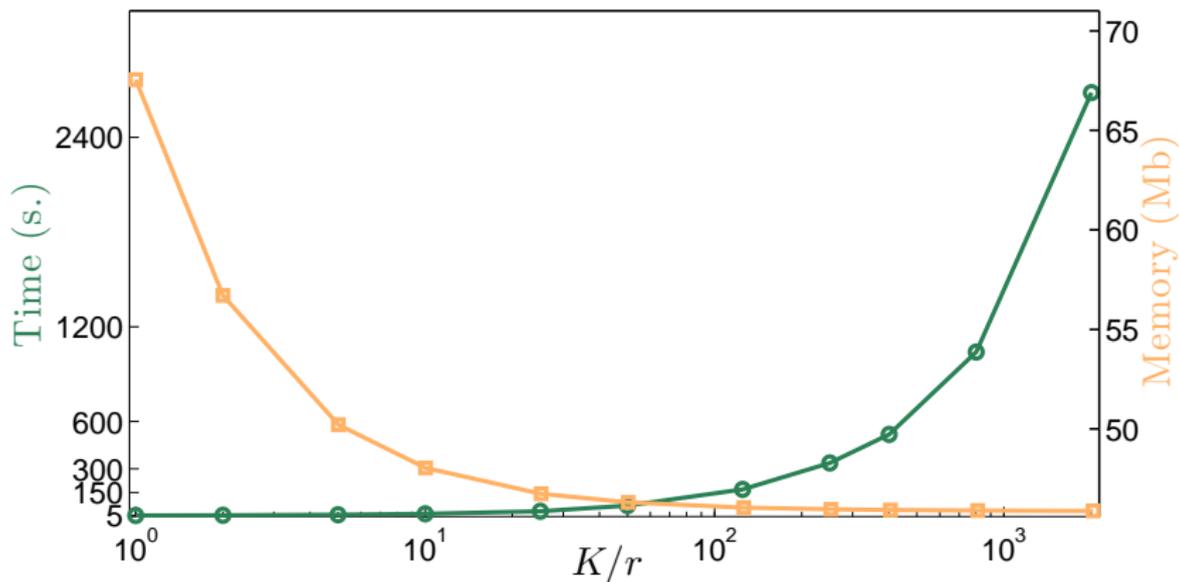
$$\text{MSE} = 8.09 \times 10^{-8}$$



Laplacian smoothing

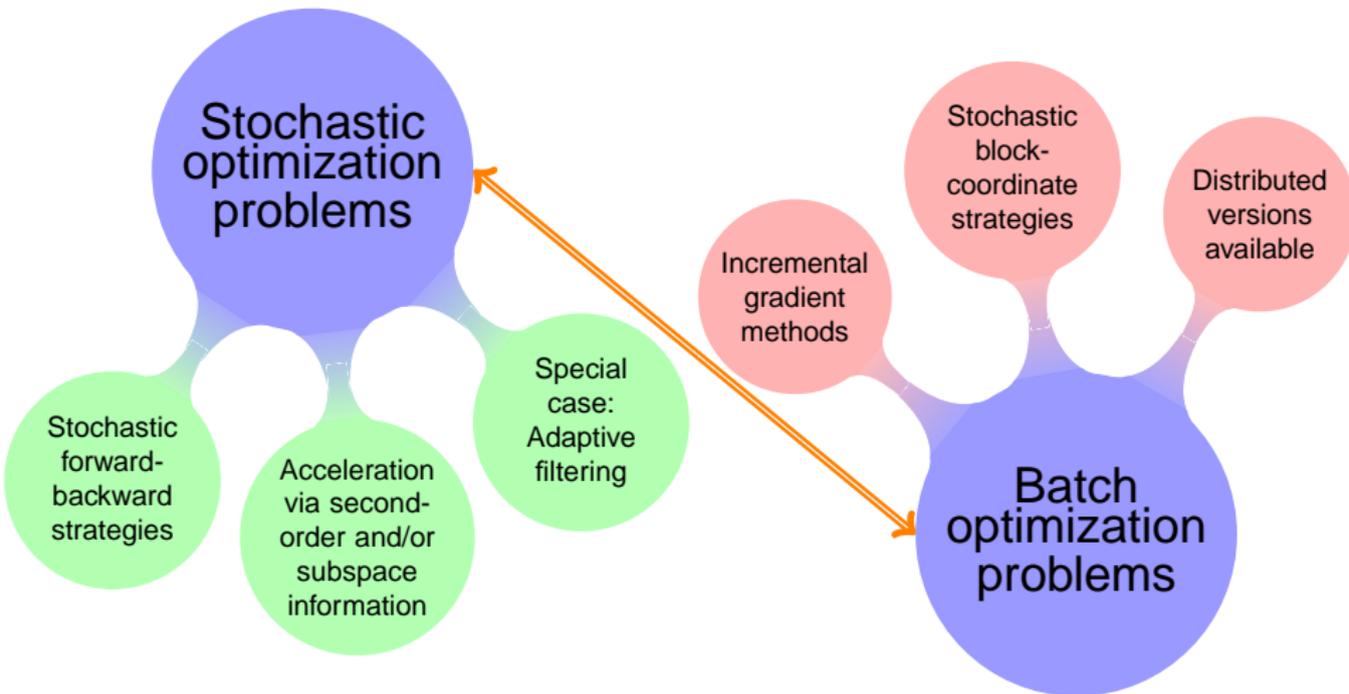
$$\text{MSE} = 5.23 \times 10^{-7}$$

- Reconstruction results using the stochastic primal-dual proximal algorithm for 3D mesh denoising from [Repetti *et al.*, 2015]:



Memory requirement, and computation time, for different number of blocks.

Conclusion





M. Pereyra, P. Schniter, E. Chouzenoux, J.-C. Pesquet, J.-Y. Tourneret, A. O. Hero and S. McLaughlin.

A Survey of Stochastic Simulation and Optimization Methods in Signal Processing.

To appear in *IEEE Journal of Selected Topics in Signal Processing.*

Available at <http://arxiv.org/abs/1505.00273>



P. Combettes and J.-C Pesquet

Stochastic Quasi-Fejér Block-Coordinate Fixed Point Iterations with Random Sweeping
SIAM Journal on Optimization, 25(2), pp. 1221-1248, 2015.



J.-C. Pesquet and A. Repetti

A Class of Randomized Primal-Dual Algorithms for Distributed Optimization
to appear in *Journal of Nonlinear and Convex Analysis*, 2015.



A. Repetti, E. Chouzenoux and J.-C. Pesquet

A Random Block-Coordinate Primal-Dual Proximal Algorithm with Application to 3D Mesh Denoising
IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2015), pp. 3561-3565, Brisbane, Australia, Apr. 19-24, 2015.



E. Chouzenoux, J.-C. Pesquet and A. Florescu.

A Stochastic 3MG Algorithm with Application to 2D Filter Identification
European Signal Processing Conference (EUSIPCO 2014), pp. 1587-1591, Lisboa, Portugal, 1-5 Sept. 2014.