

A STOCHASTIC MAJORIZATION-MINIMIZATION SUBSPACE ALGORITHM WITH APPLICATION TO FILTER IDENTIFICATION PROBLEMS

Émilie Chouzenoux and Jean-Christophe Pesquet

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

3 June 2016



Introduction

STOCHASTIC PROBLEM

$$\underset{\mathbf{h} \in \mathbb{R}^N}{\text{minimize}} \quad \frac{1}{2} \mathbb{E}(\|\mathbf{y}_n - \mathbf{X}_n^\top \mathbf{h}\|^2) + \Psi(\mathbf{h})$$

$(\mathbf{X}_n)_{n \geq 1}$ random matrices in $\mathbb{R}^{N \times Q}$,
 $(\mathbf{y}_n)_{n \geq 1}$ random vectors in \mathbb{R}^Q with

$$(\forall n \in \mathbb{N}^*) \quad \mathbb{E}(\|\mathbf{y}_n\|^2) = \varrho$$

$$\mathbb{E}(\mathbf{X}_n \mathbf{y}_n) = \mathbf{r}$$

$$\mathbb{E}(\mathbf{X}_n \mathbf{X}_n^\top) = \mathbf{R},$$

and $\Psi: \mathbb{R}^N \rightarrow \mathbb{R}$ regularization function.

Introduction

NUMEROUS EXAMPLES:

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

How to solve the problem **efficiently** when the second-order statistics of $(\mathbf{X}_n, \mathbf{y}_n)_{n \geq 1}$ are estimated **online**,
in an **adaptive** manner ?

Introduction

NUMEROUS EXAMPLES:

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

How to solve the problem **efficiently** when the second-order statistics of $(\mathbf{X}_n, \mathbf{y}_n)_{n \geq 1}$ are estimated **online**,
in an **adaptive** manner ?

⇒ Majorize-Minimize approach.

Outline

- * PROBLEM FORMULATION
 - ▶ Stochastic approximation of the criterion
 - ▶ Form of the regularization function
- * BATCH MAJORIZATION-MINIMIZATION SUBSPACE ALGORITHM
 - ▶ Quadratic tangent majorant
 - ▶ Subspace acceleration strategy
 - ▶ Convergence results
- * STOCHASTIC MAJORIZATION-MINIMIZATION SUBSPACE ALGORITHM
 - ▶ Proposed method
 - ▶ Complexity study
 - ▶ Convergence results
- * APPLICATIONS TO FILTER IDENTIFICATION
 - ▶ Online estimation of 2D blur
 - ▶ Sparse adaptive filtering

Problem formulation

Form of the objective function

Objective function at iteration $n \in \mathbb{N}^*$:

$$(\forall h \in \mathbb{R}^N) \quad F_n(h) = \frac{1}{2} \rho_n - r_n^\top h + \frac{1}{2} h^\top R_n h + \Psi(h)$$

* **Batch case:**

$$\rho_n \equiv \varrho, r_n \equiv r \text{ and } R_n \equiv R.$$

* **Online case:**

$$\rho_n = \frac{1}{n} \sum_{k=1}^n \|y_k\|^2, r_n = \frac{1}{n} \sum_{k=1}^n X_k y_k, R_n = \frac{1}{n} \sum_{k=1}^n X_k X_k^\top.$$

Form of the regularization function

$$(\forall \mathbf{h} \in \mathbb{R}^N) \quad \Psi(\mathbf{h}) = \underbrace{\frac{1}{2} \mathbf{h}^\top \mathbf{V}_0 \mathbf{h} - \mathbf{v}_0^\top \mathbf{h}}_{\text{elastic net penalization}} + \sum_{s=1}^S \psi_s(\|\mathbf{V}_s \mathbf{h} - \mathbf{v}_s\|)$$

where $\mathbf{v}_0 \in \mathbb{R}^N$, $\mathbf{V}_0 \in \mathbb{R}^{N \times N}$ symmetric positive semi-definite,
for every $s \in \{1, \dots, S\}$, $\mathbf{v}_s \in \mathbb{R}^{P_s}$, $\mathbf{V}_s \in \mathbb{R}^{P_s \times N}$, $\psi_s: \mathbb{R} \rightarrow \mathbb{R}$

~~ ability to take into account linear operators .

Form of the regularization function

$$(\forall \mathbf{h} \in \mathbb{R}^N) \quad \Psi(\mathbf{h}) = \underbrace{\frac{1}{2} \mathbf{h}^\top \mathbf{V}_0 \mathbf{h} - \mathbf{v}_0^\top \mathbf{h}}_{\text{elastic net penalization}} + \sum_{s=1}^S \psi_s(\|\mathbf{V}_s \mathbf{h} - \mathbf{v}_s\|)$$

where $\mathbf{v}_0 \in \mathbb{R}^N$, $\mathbf{V}_0 \in \mathbb{R}^{N \times N}$ symmetric positive semi-definite,
for every $s \in \{1, \dots, S\}$, $\mathbf{v}_s \in \mathbb{R}^{P_s}$, $\mathbf{V}_s \in \mathbb{R}^{P_s \times N}$, $\psi_s: \mathbb{R} \rightarrow \mathbb{R}$

~~ ability to take into account linear operators .

Assumptions on $(\psi_s)_{1 \leq s \leq S}$:

- (i) For every $s \in \{1, \dots, S\}$, ψ_s is an even lower-bounded function, which is continuously differentiable, and $\lim_{\substack{t \rightarrow 0 \\ t \neq 0}} \dot{\psi}_s(t)/t \in \mathbb{R}$,

where $\dot{\psi}_s$ denotes the derivative of ψ_s .

- (ii) For every $s \in \{1, \dots, S\}$, $\psi_s(\sqrt{\cdot})$ is concave on $[0, +\infty[$.

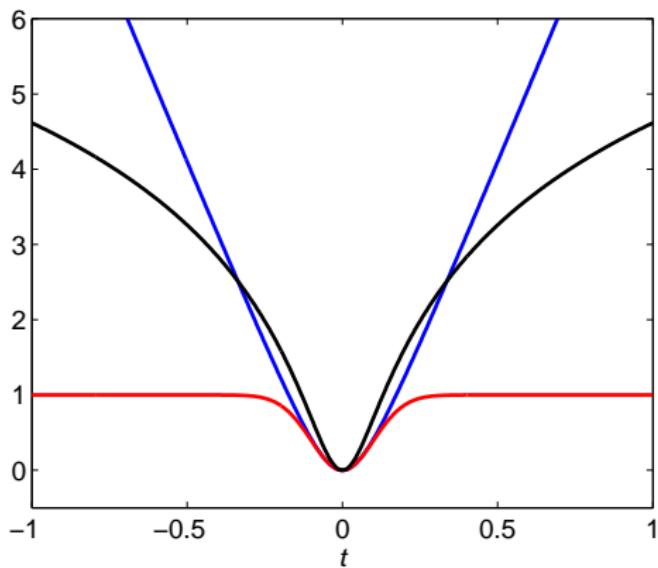
- (iii) There exists $\bar{\nu} \in [0, +\infty[$ such that $(\forall s \in \{1, \dots, S\}) (\forall t \in [0, +\infty[)$
 $0 \leq \nu_s(t) \leq \bar{\nu}$, where $\nu_s(t) = \dot{\psi}_s(t)/t$.

Examples of functions $(\psi_s)_{1 \leq s \leq S}$

	$\lambda_s^{-1} \psi_s(t)$	Type	Name
Convex	$ t - \delta_s \log(t /\delta_s + 1)$	$\ell_2 - \ell_1$	
	$\begin{cases} t^2 & \text{if } t < \delta_s \\ 2\delta_s t - \delta_s^2 & \text{otherwise} \end{cases}$	$\ell_2 - \ell_1$	Huber
	$\log(\cosh(t))$	$\ell_2 - \ell_1$	Green
	$(1 + t^2/\delta_s^2)^{\kappa_s/2} - 1$	$\ell_2 - \ell_{\kappa_s}$	
Nonconvex	$1 - \exp(-t^2/(2\delta_s^2))$	$\ell_2 - \ell_0$	Welsch
	$t^2/(2\delta_s^2 + t^2)$	$\ell_2 - \ell_0$	Geman -McClure
	$\begin{cases} 1 - (1 - t^2/(6\delta_s^2))^3 & \text{if } t \leq \sqrt{6}\delta_s \\ 1 & \text{otherwise} \end{cases}$	$\ell_2 - \ell_0$	Tukey biweight
	$\tanh(t^2/(2\delta_s^2))$	$\ell_2 - \ell_0$	Hyperbolic tangent
	$\log(1 + t^2/\delta_s^2)$	$\ell_2 - \log$	Cauchy
	$1 - \exp(1 - (1 + t^2/(2\delta_s^2))^{\kappa_s/2})$	$\ell_2 - \ell_{\kappa_s} - \ell_0$	Chouzenoux

$$(\lambda_s, \delta_s) \in]0, +\infty[^2, \kappa_s \in [1, 2]$$

Examples of functions $(\psi_s)_{1 \leq s \leq S}$



$$\psi_s(t) = (1 + \frac{t^2}{\delta^2})^{1/2} - 1, \quad \psi_s(t) = \log\left(1 + \frac{t^2}{\delta^2}\right), \quad \psi_s(t) = 1 - \exp(-\frac{t^2}{2\delta^2}).$$

Batch majorize-minimize subspace algorithm

Majorize-Minimize principle

1. Find a tractable surrogate for $F \rightsquigarrow$ Majorization step

\rightsquigarrow Quadratic tangent majorant

$$\begin{aligned} (\forall \mathbf{h} \in \mathbb{R}^N) \quad \Theta(\mathbf{h}, \mathbf{h}_n) &= F(\mathbf{h}_n) + \nabla F(\mathbf{h}_n)^\top (\mathbf{h} - \mathbf{h}_n) \\ &\quad + \frac{1}{2} (\mathbf{h} - \mathbf{h}_n)^\top \mathbf{A}(\mathbf{h}_n) (\mathbf{h} - \mathbf{h}_n), \\ &\geq F(\mathbf{h}), \end{aligned}$$

where $\mathbf{A}(\mathbf{h}) = \mathbf{R} + \mathbf{V}_0 + \mathbf{V}^\top \text{Diag}(\mathbf{b}(\mathbf{h})) \mathbf{V} \in \mathbb{R}^{N \times N}$

with $\mathbf{V} = [\mathbf{V}_1^\top \dots \mathbf{V}_S^\top]^\top \in \mathbb{R}^{P \times N}$, $P = P_1 + \dots + P_S$

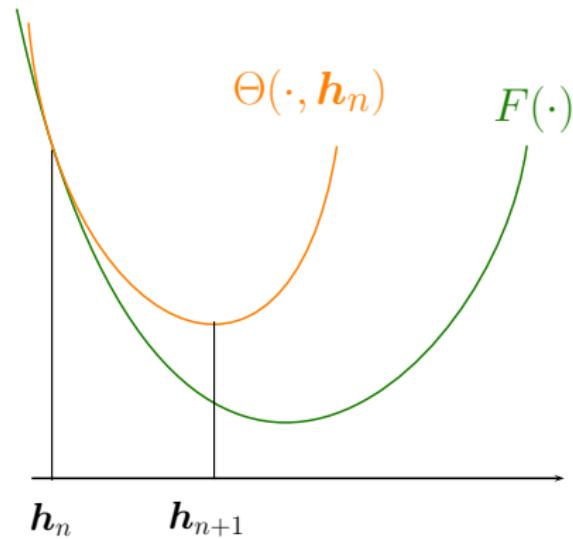
and $\mathbf{b}(\mathbf{h}) = (b_i(\mathbf{h}))_{1 \leq i \leq P} \in \mathbb{R}^P$ such that

$$(\forall s \in \{1, \dots, S\})(\forall p \in \{1, \dots, P_s\})$$

$$b_{P_1 + \dots + P_{s-1} + p}(\mathbf{h}) = \nu_s(\|\mathbf{V}_s \mathbf{h} - \mathbf{v}_s\|).$$

Majorize-Minimize principle

1. Find a tractable surrogate for $F \rightsquigarrow$ Majorization step



Subspace acceleration

2. Minimize in a subspace \rightsquigarrow Minimization step

$$(\forall n \in \mathbb{N}^*) \quad \mathbf{h}_{n+1} \in \underset{\mathbf{h} \in \text{span } \mathcal{D}_n}{\text{Argmin}} \Theta(\mathbf{h}, \mathbf{h}_n),$$

with $\mathcal{D}_n \in \mathbb{R}^{N \times M_n}$.

- ▶ $\text{rank}(\mathcal{D}_n) = N \Rightarrow$ half-quadratic algorithm
- ▶ M_n small \Rightarrow low-complexity per iteration.

Typical choice:

$$\mathcal{D}_n = \begin{cases} [-\nabla F(\mathbf{h}_n), \mathbf{h}_n, \mathbf{h}_n - \mathbf{h}_{n-1}] & \text{if } n > 1 \\ [-\nabla F(\mathbf{h}_1), \mathbf{h}_1] & \text{if } n = 1 \end{cases}$$

\rightsquigarrow **3MG** algorithm

(similar ideas in TWIST, FISTA, NLCG, L-BFGS, ...)

Batch MM subspace algorithm

Initialize $D_0, u_0, h_1 = D_0 u_0$

$$D_0^R = R D_0, D_0^{V_0} = V_0 D_0, D_0^V = V D_0$$

For $n = 1, 2, \dots$

$$c(h_n) = r + v_0 + V^\top \text{Diag}(b(h)) v$$

$$D_{n-1}^A = D_{n-1}^R + D_{n-1}^{V_0} + V^\top \text{Diag}(b(h_n)) D_{n-1}^V$$

$$\nabla F(h_n) = D_{n-1}^A u_{n-1} - c(h_n)$$

Set D_n using $\nabla F(h_n)$

$$D_n^R = R D_n, D_n^{V_0} = V_0 D_n, D_n^V = V D_n$$

$$B_n = D_n^\top (D_n^{V_0} + D_n^R) + (D_n^V)^\top \text{Diag}(b(h_n)) D_n^V$$

$$u_n = B_n^\dagger D_n^\top c(h_n)$$

$$h_{n+1} = D_n u_n$$

Convergence theorem

Let assume that:

1. For every $n \in \mathbb{N}^*$, $\{\nabla F(\mathbf{h}_n), \mathbf{h}_n\} \subset \text{span } \mathbf{D}_n$,
2. $\mathbf{R} + \mathbf{V}_0$ is a positive definite matrix.

Then, the following hold:

- $\|\nabla F(\mathbf{h}_n)\| \rightarrow 0$ and $F(\mathbf{h}_n) \searrow F(\hat{\mathbf{h}})$ where $\hat{\mathbf{h}}$ is a critical point of F .
- If the functions $(\psi_s)_{1 \leq s \leq S}$ are convex, then $(\mathbf{h}_n)_{n \geq 1}$ converges to the unique (global) minimizer $\hat{\mathbf{h}}$ of F
- If F satisfies the Kurdyka-Łojasiewicz inequality, then the sequence $(\mathbf{h}_n)_{n \geq 1}$ converges to a critical point of F .

Stochastic majorize-minimize subspace algorithm

Stochastic approximation of the criterion

Estimate of the objective function at iteration $n \in \mathbb{N}^*$:

$$\begin{aligned} (\forall h \in \mathbb{R}^N) \quad F_n(h) &= \frac{1}{2n} \sum_{k=1}^n \|y_k - X_k^\top h\|^2 + \Psi(h) \\ &= \frac{1}{2} \rho_n - r_n^\top h + \frac{1}{2} h^\top R_n h + \Psi(h) \end{aligned}$$

with $\rho_n = \frac{1}{n} \sum_{k=1}^n \|y_k\|^2$, $r_n = \frac{1}{n} \sum_{k=1}^n X_k y_k$, and

$$R_n = \frac{1}{n} \sum_{k=1}^n X_k X_k^\top.$$

Stochastic approximation of the criterion

Estimate of the objective function at iteration $n \in \mathbb{N}^*$:

$$\begin{aligned} (\forall h \in \mathbb{R}^N) \quad F_n(h) &= \frac{1}{2n} \sum_{k=1}^n \|y_k - X_k^\top h\|^2 + \Psi(h) \\ &= \frac{1}{2} \rho_n - r_n^\top h + \frac{1}{2} h^\top R_n h + \Psi(h) \end{aligned}$$

with $\rho_n = \frac{1}{n} \sum_{k=1}^n \|y_k\|^2$, $r_n = \frac{1}{n} \sum_{k=1}^n X_k y_k$, and

$$R_n = \frac{1}{n} \sum_{k=1}^n X_k X_k^\top.$$

- How to make the method adaptive to changes in the input statistics?

Stochastic approximation of the criterion

Estimate of the objective function at iteration $n \in \mathbb{N}^*$:

$$\begin{aligned} (\forall \mathbf{h} \in \mathbb{R}^N) \quad F_n(\mathbf{h}) &= \frac{1}{2\bar{\vartheta}_n} \sum_{k=1}^n \vartheta^{n-k} \|\mathbf{y}_k - \mathbf{X}_k^\top \mathbf{h}\|^2 + \Psi(\mathbf{h}) \\ &= \frac{1}{2} \rho_n - \mathbf{r}_n^\top \mathbf{h} + \frac{1}{2} \mathbf{h}^\top \mathbf{R}_n \mathbf{h} + \Psi(\mathbf{h}) \end{aligned}$$

with $\rho_n = \frac{1}{\bar{\vartheta}_n} \sum_{k=1}^n \vartheta^{n-k} \|\mathbf{y}_k\|^2$, $\mathbf{r}_n = \frac{1}{\bar{\vartheta}_n} \sum_{k=1}^n \vartheta^{n-k} \mathbf{X}_k \mathbf{y}_k$, and

$$\mathbf{R}_n = \frac{1}{\bar{\vartheta}_n} \sum_{k=1}^n \vartheta^{n-k} \mathbf{X}_k \mathbf{X}_k^\top.$$

~~~ **Forgetting factor:**

$$\bar{\vartheta}_n = \sum_{k=0}^{n-1} \vartheta^k = \begin{cases} n & \text{if } \vartheta = 1 \\ \frac{1 - \vartheta^n}{1 - \vartheta} & \text{if } \vartheta \in ]0, 1[ \end{cases}, \quad \text{and} \quad \vartheta \in ]0, 1].$$

## Stochastic MM subspace algorithm

Estimate of the objective function at iteration  $n \in \mathbb{N}^*$ :

$$(\forall \mathbf{h} \in \mathbb{R}^N) \quad F_n(\mathbf{h}) = \frac{1}{2\bar{\vartheta}_n} \sum_{k=1}^n \vartheta^{n-k} \|\mathbf{y}_k - \mathbf{X}_k^\top \mathbf{h}\|^2 + \Psi(\mathbf{h}).$$

1. Find a tractable surrogate for  $F_n$   
~~~ Quadratic tangent majorant

$$\begin{aligned} (\forall \mathbf{h} \in \mathbb{R}^N) \quad \Theta_n(\mathbf{h}, \mathbf{h}_n) &= F_n(\mathbf{h}_n) + \nabla F_n(\mathbf{h}_n)^\top (\mathbf{h} - \mathbf{h}_n) \\ &\quad + \frac{1}{2} (\mathbf{h} - \mathbf{h}_n)^\top \mathbf{A}_n(\mathbf{h}_n) (\mathbf{h} - \mathbf{h}_n), \\ &\geq F_n(\mathbf{h}), \end{aligned}$$

where $\mathbf{A}_n(\mathbf{h}) = \mathbf{R}_n + \mathbf{V}_0 + \mathbf{V}^\top \text{Diag}(\mathbf{b}(\mathbf{h})) \mathbf{V} \in \mathbb{R}^{N \times N}$.

Stochastic MM subspace algorithm

Estimate of the objective function at iteration $n \in \mathbb{N}^*$:

$$(\forall \mathbf{h} \in \mathbb{R}^N) \quad F_n(\mathbf{h}) = \frac{1}{2\bar{\vartheta}_n} \sum_{k=1}^n \vartheta^{n-k} \|\mathbf{y}_k - \mathbf{X}_k^\top \mathbf{h}\|^2 + \Psi(\mathbf{h}).$$

1. Find a tractable surrogate for F_n
2. Minimize in a subspace

$$(\forall n \in \mathbb{N}^*) \quad \mathbf{h}_{n+1} \in \underset{\mathbf{h} \in \text{span } \mathbf{D}_n}{\operatorname{Argmin}} \Theta_n(\mathbf{h}, \mathbf{h}_n),$$

with $\mathbf{D}_n \in \mathbb{R}^{N \times M_n}$.

~~~ Stochastic 3MG algorithm :

$$\mathbf{D}_n = \begin{cases} [-\nabla F_n(\mathbf{h}_n), \mathbf{h}_n, \mathbf{h}_n - \mathbf{h}_{n-1}] & \text{if } n > 1 \\ [-\nabla F_n(\mathbf{h}_1), \mathbf{h}_1] & \text{if } n = 1 \end{cases}$$

## Stochastic MM subspace algorithm

Estimate of the objective function at iteration  $n \in \mathbb{N}^*$ :

$$(\forall h \in \mathbb{R}^N) \quad F_n(h) = \frac{1}{2\bar{\vartheta}_n} \sum_{k=1}^n \vartheta^{n-k} \|y_k - X_k^\top h\|^2 + \Psi(h).$$

1. Find a tractable surrogate for  $F_n$
2. Minimize in a subspace
3. Perform recursive updates of the second-order statistics

$$(\forall n \in \mathbb{N}^*) \quad \mathbf{r}_n = \mathbf{r}_{n-1} + \frac{1}{\bar{\vartheta}_n} (\mathbf{X}_n \mathbf{y}_n - \mathbf{r}_{n-1})$$

$$\mathbf{R}_n = \mathbf{R}_{n-1} + \frac{1}{\bar{\vartheta}_n} (\mathbf{X}_n \mathbf{X}_n^\top - \mathbf{R}_{n-1}).$$

## Stochastic MM subspace algorithm

$$\mathbf{r}_0 = \mathbf{0}, \mathbf{R}_0 = \mathbf{0}$$

Initialize  $\mathbf{D}_0, \mathbf{u}_0$

$$\mathbf{h}_1 = \mathbf{D}_0 \mathbf{u}_0, \mathbf{D}_0^{\mathbf{R}} = \mathbf{0}, \mathbf{D}_0^{V_0} = V_0 \mathbf{D}_0, \mathbf{D}_0^V = V \mathbf{D}_0$$

For all  $n = 1, \dots$

$$\mathbf{r}_n = \mathbf{r}_{n-1} + \frac{1}{\vartheta_n} (\mathbf{X}_n \mathbf{y}_n - \mathbf{r}_{n-1})$$

$$\mathbf{c}_n(\mathbf{h}_n) = \mathbf{r}_n + \mathbf{v}_0 + V^\top \text{Diag}(\mathbf{b}(\mathbf{h}_n)) \mathbf{v}$$

$$\begin{aligned} \mathbf{D}_{n-1}^{\mathbf{A}} = & (1 - \frac{1}{\vartheta_n}) \mathbf{D}_{n-1}^{\mathbf{R}} + \frac{1}{\vartheta_n} \mathbf{X}_n (\mathbf{X}_n^\top \mathbf{D}_{n-1}) \\ & + \mathbf{D}_{n-1}^{V_0} + V^\top \text{Diag}(\mathbf{b}(\mathbf{h}_n)) \mathbf{D}_{n-1}^V \end{aligned}$$

$$\nabla F_n(\mathbf{h}_n) = \mathbf{D}_{n-1}^{\mathbf{A}} \mathbf{u}_{n-1} - \mathbf{c}_n(\mathbf{h}_n)$$

$$\mathbf{R}_n = \mathbf{R}_{n-1} + \frac{1}{\vartheta_n} (\mathbf{X}_n \mathbf{X}_n^\top - \mathbf{R}_{n-1})$$

Set  $\mathbf{D}_n$  using  $\nabla F_n(\mathbf{h}_n)$

$$\mathbf{D}_n^{\mathbf{R}} = \mathbf{R}_n \mathbf{D}_n, \mathbf{D}_n^{V_0} = V_0 \mathbf{D}_n, \mathbf{D}_n^V = V \mathbf{D}_n$$

$$\mathbf{B}_n = \mathbf{D}_n^\top (\mathbf{D}_n^{\mathbf{R}} + \mathbf{D}_n^{V_0}) + (\mathbf{D}_n^V)^\top \text{Diag}(\mathbf{b}(\mathbf{h}_n)) \mathbf{D}_n^V$$

$$\mathbf{u}_n = \mathbf{B}_n^\dagger \mathbf{D}_n^\top \mathbf{c}_n(\mathbf{h}_n)$$

$$\mathbf{h}_{n+1} = \mathbf{D}_n \mathbf{u}_n$$

## Computational complexity

- ▶ If  $V = I_N$ , number of multiplications per iteration:

$$N^2(4M_n + Q)/2,$$

where  $N$  length of  $h$ ,  $Q$  column dimension of  $\mathbf{X}_n$ ,  $M_n$  subspace dimension, and  $N \gg \max\{M_n, M_{n-1}, Q\}$ .

## Computational complexity

- ▶ If  $\mathbf{V} = \mathbf{I}_N$ , number of multiplications per iteration:

$$N^2(4M_n + Q)/2,$$

where  $N$  length of  $\mathbf{h}$ ,  $Q$  column dimension of  $\mathbf{X}_n$ ,  $M_n$  subspace dimension, and  $N \gg \max\{M_n, M_{n-1}, Q\}$ .

- ▶ If  $\mathbf{V} \neq \mathbf{I}_N$ , number of multiplications per iteration:

$$N \left( \underbrace{P(M_n + M_{n-1} + 1)}_{\text{upper bound on the complexity induced by linear operators}} + N(4M_n + Q)/2 \right),$$

where  $P$  row dimension of  $\mathbf{V}$ .

## Computational complexity

- ▶ If  $\mathbf{V} = \mathbf{I}_N$ , number of multiplications per iteration:

$$N^2(4M_n + Q)/2,$$

where  $N$  length of  $\mathbf{h}$ ,  $Q$  column dimension of  $\mathbf{X}_n$ ,  $M_n$  subspace dimension, and  $N \gg \max\{M_n, M_{n-1}, Q\}$ .

- ▶ If  $\mathbf{V} \neq \mathbf{I}_N$ , number of multiplications per iteration:

$$N \left( \underbrace{P(M_n + M_{n-1} + 1)}_{\text{upper bound on the complexity induced by linear operators}} + N(4M_n + Q)/2 \right),$$

where  $P$  row dimension of  $\mathbf{V}$ .

- ▶ Further complexity reduction possible by taking into account the structure of  $\mathbf{D}_n$ ,  
e.g. 3MG,  $\mathbf{V} = \mathbf{I}_N$ ,  $Q = 1 \equiv$  complexity of RLS.

## Convergence results: Assumptions

Let  $(\Omega, \mathcal{F}, P)$  be the underlying probability space. For every  $n \in \mathbb{N}^*$ , let  $\mathcal{X}_n = \sigma((\mathbf{X}_k, \mathbf{y}_k)_{1 \leq k \leq n})$  be the sub-sigma algebra of  $\mathcal{F}$  generated by  $(\mathbf{X}_k, \mathbf{y}_k)_{1 \leq k \leq n}$ .

1.  $\mathbf{R} + \mathbf{V}_0$  is a positive definite matrix.
2.  $((\mathbf{X}_n, \mathbf{y}_n))_{n \geq 1}$  is a stationary ergodic sequence and, for every  $n \in \mathbb{N}^*$ , the elements of  $\mathbf{X}_n$  and the components of  $\mathbf{y}_n$  have finite fourth-order moments.
3.  $(\forall n \in \mathbb{N}^*) E(\|\mathbf{y}_{n+1}\|^2 | \mathcal{X}_n) = \varrho$ ,  $E(\mathbf{X}_{n+1}\mathbf{y}_{n+1} | \mathcal{X}_n) = \mathbf{r}$  and  $E(\mathbf{X}_{n+1}\mathbf{X}_{n+1}^\top | \mathcal{X}_n) = \mathbf{R}$ .
4. For every  $n \in \mathbb{N}^*$ ,  $\{\nabla F_n(\mathbf{h}_n), \mathbf{h}_n\} \subset \text{span } \mathbf{D}_n$ .
5.  $\mathbf{h}_1$  is  $\mathcal{X}_1$ -measurable and, for every  $n \in \mathbb{N}^*$ ,  $\mathbf{D}_n$  is  $\mathcal{X}_n$ -measurable.

## Convergence results: Theorem

- $(\rho_n)_{n \geq 1}$ ,  $(\mathbf{R}_n)_{n \geq 1}$ , and  $(\mathbf{r}_n)_{n \geq 1}$  converge almost surely to  $\varrho$ ,  $\mathbf{R}$  and  $\mathbf{r}$ , respectively.
- The set of cluster points of  $(\mathbf{h}_n)_{n \geq 1}$  is almost surely a nonempty compact connected set.
- Any element of this set is almost surely a critical point of  $F$ .
- If the functions  $(\psi_s)_{1 \leq s \leq S}$  are convex, then the sequence  $(\mathbf{h}_n)_{n \geq 1}$  converges almost surely to the unique (global) minimizer of  $F$ .

## Convergence results : Rate analysis

Assume that  $\Psi$  is convex and twice differentiable.

Let  $\epsilon \in ]0, +\infty[$  be such that  $\epsilon I_N \prec R + V_0$ .

- There exists P-a.s.  $n_\epsilon \in \mathbb{N}^*$  such that, for every  $n \geq n_\epsilon$ ,

$$F_n(\mathbf{h}_{n+1}) - \inf F_n \leq \theta_n (F_n(\mathbf{h}_n) - \inf F_n)$$

where  $\theta_n = 1 - (1 + \epsilon)^{-1} \tilde{\theta}_n$ , with

$$\tilde{\theta}_n = \frac{(\nabla F_n(\mathbf{h}_n))^\top \mathbf{D}_n (\mathbf{D}_n^\top \mathbf{A}_n(\mathbf{h}_n) \mathbf{D}_n)^\dagger \mathbf{D}_n^\top \nabla F_n(\mathbf{h}_n)}{(\nabla F_n(\mathbf{h}_n))^\top (\nabla^2 F_n(\mathbf{h}_n))^{-1} \nabla F_n(\mathbf{h}_n)}.$$

## Convergence results : Rate analysis

Assume that  $\Psi$  is convex and twice differentiable.

Let  $\epsilon \in ]0, +\infty[$  be such that  $\epsilon I_N \prec R + V_0$ .

- There exists P-a.s.  $n_\epsilon \in \mathbb{N}^*$  such that, for every  $n \geq n_\epsilon$ ,

$$F_n(\mathbf{h}_{n+1}) - \inf F_n \leq \theta_n (F_n(\mathbf{h}_n) - \inf F_n)$$

- $\theta_n \in [\underline{\theta}_n, \bar{\theta}_n]$  with

$$\begin{cases} \underline{\theta}_n = 1 - (1 + \epsilon)^{-1} \underline{\kappa}_n^{-1} > 0, \\ \bar{\theta}_n = 1 - (1 + \epsilon)^{-1} \bar{\kappa}_n^{-1} \left( 1 - \left( \frac{\bar{\sigma}_n - \underline{\sigma}_n}{\bar{\sigma}_n + \underline{\sigma}_n} \right)^2 \right) < 1. \end{cases}$$

where  $\underline{\kappa}_n$  (resp.  $\bar{\kappa}_n$ ) is the minimum (resp. maximum) eigenvalue of  $(\mathbf{A}_n(\mathbf{h}_n))^{\frac{1}{2}} (\nabla^2 F_n(\mathbf{h}_n))^{-1} (\mathbf{A}_n(\mathbf{h}_n))^{\frac{1}{2}}$  and  $\underline{\sigma}_n$  (resp.  $\bar{\sigma}_n$ ) denotes the minimum (resp. maximum) eigenvalue of the Hessian of  $F_n$  at  $\mathbf{h}_n$ .

## Convergence results : Rate analysis

Assume that  $\Psi$  is convex and twice differentiable.

Let  $\epsilon \in ]0, +\infty[$  be such that  $\epsilon I_N \prec R + V_0$ .

- There exists P-a.s.  $n_\epsilon \in \mathbb{N}^*$  such that, for every  $n \geq n_\epsilon$ ,

$$F_n(\mathbf{h}_{n+1}) - \inf F_n \leq \theta_n (F_n(\mathbf{h}_n) - \inf F_n)$$

- $\theta_n \in [\underline{\theta}_n, \bar{\theta}_n]$  with

$$\begin{cases} \underline{\theta}_n = 1 - (1 + \epsilon)^{-1} \underline{\kappa}_n^{-1} > 0, \\ \bar{\theta}_n = 1 - (1 + \epsilon)^{-1} \bar{\kappa}_n^{-1} \left( 1 - \left( \frac{\bar{\sigma}_n - \underline{\sigma}_n}{\bar{\sigma}_n + \underline{\sigma}_n} \right)^2 \right) < 1. \end{cases}$$

- ▶ slowest rate  $\bar{\theta}_n$  for gradient subspace
- ▶ fastest rate  $\underline{\theta}_n$  for full subspace.

## Convergence results : Rate analysis in the batch case

Assume that  $\Psi$  is convex and twice differentiable.

Let  $\epsilon \in ]0, +\infty[$  be such that  $\epsilon I_N \prec R + V_0$ .

- There exists  $n_\epsilon \in \mathbb{N}^*$ ,  $\mu \in \mathbb{R}^+$  and  $\vartheta \in ]0, 1[$  such that, for every  $n \geq n_\epsilon$ ,

$$F(\mathbf{h}_n) - \inf F \leq \mu \vartheta^n.$$

## Convergence results : Rate analysis in the batch case

Assume that  $\Psi$  is convex and twice differentiable.

Let  $\epsilon \in ]0, +\infty[$  be such that  $\epsilon \mathbf{I}_N \prec \mathbf{R} + \mathbf{V}_0$ .

- There exists  $n_\epsilon \in \mathbb{N}^*$ ,  $\mu \in \mathbb{R}^+$  and  $\vartheta \in ]0, 1[$  such that, for every  $n \geq n_\epsilon$ ,

$$F(\mathbf{h}_n) - \inf F \leq \mu \vartheta^n.$$

- The worst-case geometrical decay rate is

$$\vartheta = 1 - \frac{1}{(1 + \epsilon)\bar{\kappa}_{\max}} \left( 1 - \left( \frac{\bar{\eta} - \underline{\eta} + 2\epsilon}{\bar{\eta} + \underline{\eta}} \right)^2 \right).$$

where  $\underline{\eta}$  (resp.  $\bar{\eta}_n$ ) denotes the minimum (resp. maximum) eigenvalue of  $\mathbf{R} + \mathbf{V}_0$  (resp.  $\mathbf{R} + \mathbf{V}_0 + \bar{\nu} \mathbf{V}^\top \mathbf{V}$ ).

## Convergence results : Rate analysis in the batch case

Assume that  $\Psi$  is convex and twice differentiable.

Let  $\epsilon \in ]0, +\infty[$  be such that  $\epsilon \mathbf{I}_N \prec \mathbf{R} + \mathbf{V}_0$ .

- There exists  $n_\epsilon \in \mathbb{N}^*$ ,  $\mu \in \mathbb{R}^+$  and  $\vartheta \in ]0, 1[$  such that, for every  $n \geq n_\epsilon$ ,

$$F(\mathbf{h}_n) - \inf F \leq \mu \vartheta^n.$$

- The worst-case geometrical decay rate is

$$\vartheta = 1 - \frac{1}{(1 + \epsilon)\bar{\kappa}_{\max}} \left( 1 - \left( \frac{\bar{\eta} - \underline{\eta} + 2\epsilon}{\bar{\eta} + \underline{\eta}} \right)^2 \right).$$

where  $\underline{\eta}$  (resp.  $\bar{\eta}_n$ ) denotes the minimum (resp. maximum) eigenvalue of  $\mathbf{R} + \mathbf{V}_0$  (resp.  $\mathbf{R} + \mathbf{V}_0 + \bar{\nu} \mathbf{V}^\top \mathbf{V}$ ).

Linear convergence of  $(\mathbf{h}_n)_{n \in \mathbb{N}}$  in the batch case.

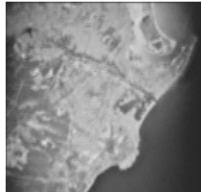
# Application to filter identification problems

## Application to 2D filter identification

## OBSERVATION MODEL

$$\mathbf{y} = S(\overline{\mathbf{h}})\mathbf{x} + \mathbf{w}$$

- ▶  $x \in \mathbb{R}^L$  large size original image ( $L = 4096^2$ ),
  - ▶  $\bar{h} \in \mathbb{R}^N$  unknown two-dimensional blur kernel ( $N = 21^2$ ),
  - ▶  $S(\bar{h})$  Hankel-block Hankel matrix such that  
 $S(\bar{h})x = X\bar{h}$ ,
  - ▶  $w \in \mathbb{R}^L$  realization of white  $\mathcal{N}(0, 0.03^2)$  noise  
(BSNR = 25.7 dB)
  - ▶  $y \in \mathbb{R}^L$  blurred and noisy image.

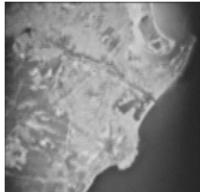


## Application to 2D filter identification

## OBSERVATION MODEL

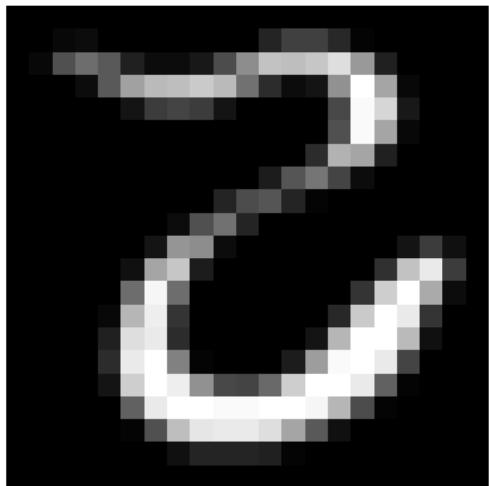
$$\mathbf{y} = S(\bar{\mathbf{h}})\mathbf{x} + \mathbf{w}$$

- ▶  $x \in \mathbb{R}^L$  large size original image ( $L = 4096^2$ ),
  - ▶  $\bar{h} \in \mathbb{R}^N$  unknown two-dimensional blur kernel ( $N = 21^2$ ),
  - ▶  $S(\bar{h})$  Hankel-block Hankel matrix such that  
 $S(\bar{h})x = X\bar{h}$ ,
  - ▶  $w \in \mathbb{R}^L$  realization of white  $\mathcal{N}(0, 0.03^2)$  noise  
(BSNR = 25.7 dB)
  - ▶  $y \in \mathbb{R}^L$  blurred and noisy image.

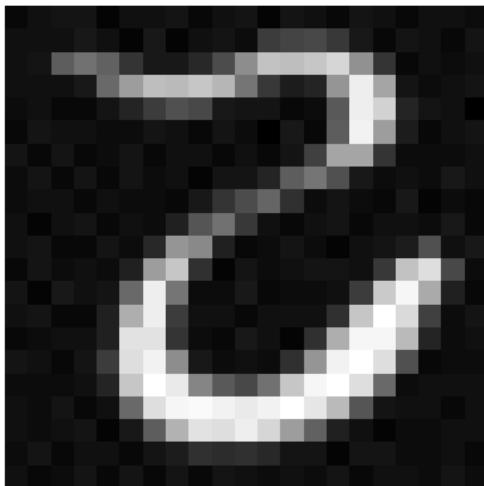


⇒ **Minimization of a penalized MSE criterion:**  $\mathbf{y}_n \in \mathbb{R}^Q$  and  $\mathbf{X}_n^\top \in \mathbb{R}^{Q \times N}$ :  $Q$  lines of  $\mathbf{y}$  and  $\mathbf{X}$ ,  $\vartheta = 1$ , and  $\Psi$  isotropic penalization on the gradient of  $\mathbf{h}$ , i.e.  $S = N$ ,  $(\forall s \in \{1, \dots, S\}) P_s = 2$ ,  $\psi_s: u \mapsto \lambda \sqrt{1 + u^2/\delta^2}$ ,  $(\lambda, \delta) \in ]0, +\infty[^2$ .

## Application to 2D filter identification



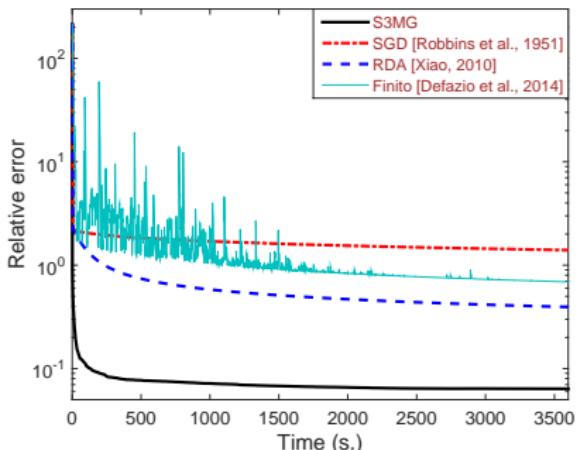
## Original blur kernel.



Estimated blur kernel, relative error 0.064.

- ▶ The regularization parameters are optimized manually.

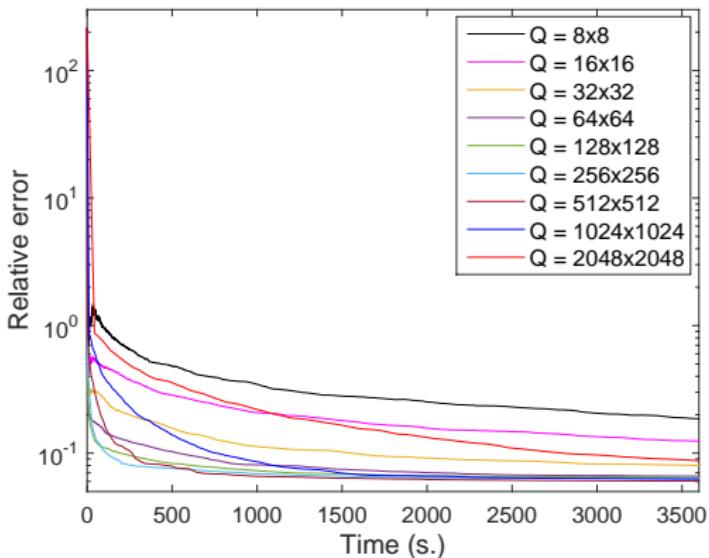
## Application to 2D filter identification



Comparison of stochastic 3MG algorithm, SGD algorithm with decreasing stepsize  $\propto n^{-1/2}$ , and SAGA/RDA algorithms with constant stepsizes.

- ▶ The stepsize values in SGD/SAGA/RDA methods are optimized manually .
- ▶ The S3MG algorithm leads to a faster convergence .

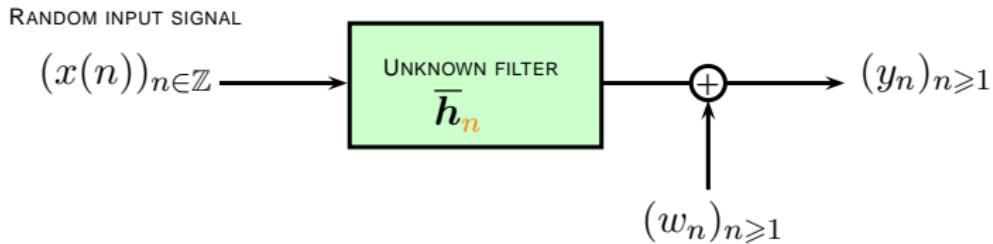
## Application to 2D filter identification



Effect of the block size  $Q$  on the convergence speed of S3MG.

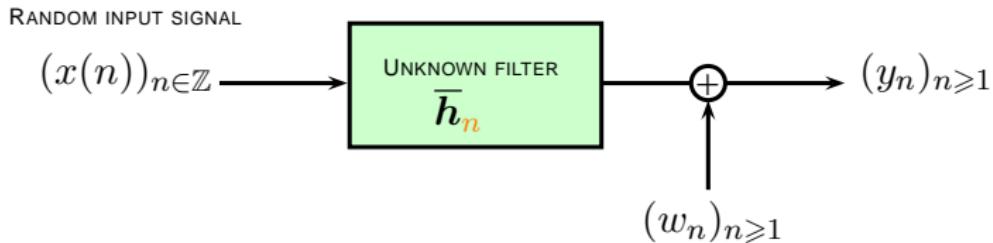
- ▶ The best trade-off is obtained for  $Q = 256 \times 256$ .

## Application to sparse adaptive filtering



$(\forall n \in \mathbb{Z}) \quad \bar{h}_n \in \mathbb{R}^N$  **sparse** filter,  $y_n \in \mathbb{R}$ ,  $w_n \in \mathbb{R}$ .

## Application to sparse adaptive filtering



$(\forall n \in \mathbb{Z}) \quad \bar{h}_n \in \mathbb{R}^N$  **sparse filter**,  $y_n \in \mathbb{R}$ ,  $w_n \in \mathbb{R}$ .

⇒ **Minimization of a penalized MSE criterion:**

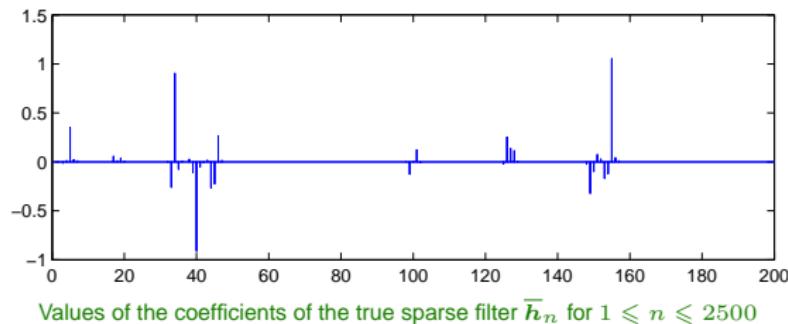
- ~~  $\mathbf{X}_n = [x(n-N+1), \dots, x(n)]^\top \in \mathbb{R}^N$ ;
- ~~ Smoothed  $\ell_0$  regularization with  $S = N$ ,  $\mathbf{v}_0 = \mathbf{0}$ ,  $\mathbf{V}_0 = \mathbf{O}_N$ , and  $(\forall s \in \{1, \dots, N\}) P_s = 1$ ,  $\mathbf{v}_s = 0$ ,  $\mathbf{V}_s \in \mathbb{R}^{1 \times N}$  the  $s$ -th vector of the canonical basis of  $\mathbb{R}^N$  and  $\psi_s : u \mapsto \lambda(1 - \exp(-u^2/(2\delta^2)))$ .

## Simulation results

$(\bar{h}_n)$  : Time-variant linear system with 200 sparse coefficients,

$(x(n))_n$  : Input sequence of 5000 random independent variables uniformly distributed on  $\{-1, +1\}$ ,

$(w_n)_n$  : White Gaussian noise with zero mean and variance 0.05.

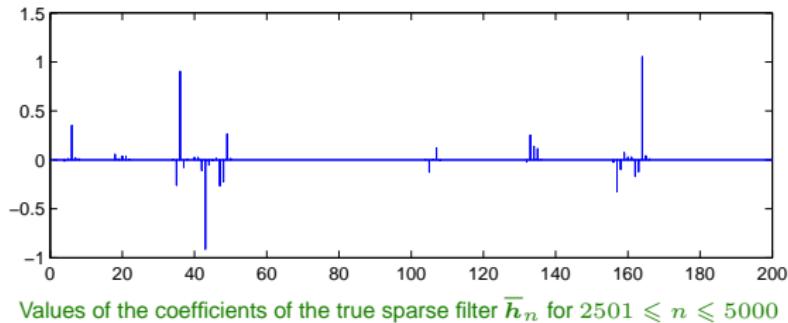


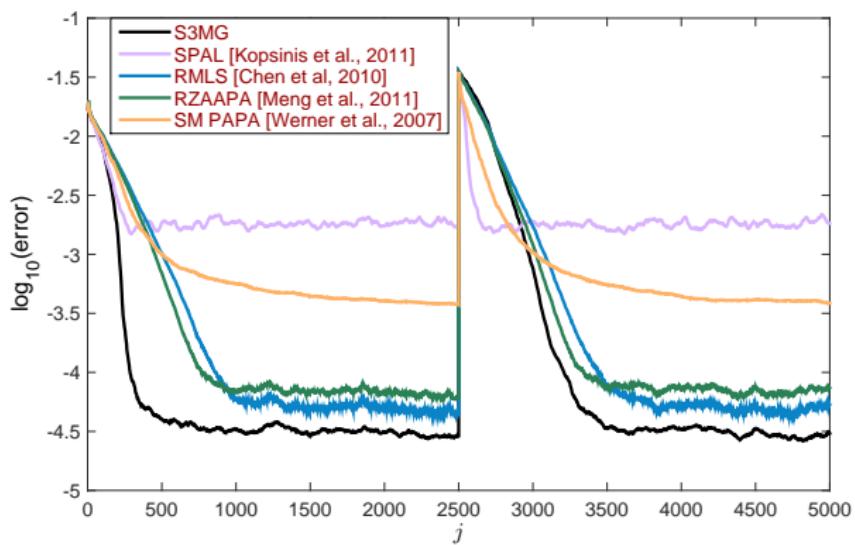
## Simulation results

( $\bar{h}_n$ ) : Time-variant linear system with 200 sparse coefficients,

$(x(n))_n$  : Input sequence of 5000 random independent variables uniformly distributed on  $\{-1, +1\}$ ,

$(w_n)_n$  : White Gaussian noise with zero mean and variance 0.05.





Estimation error along time, for various sparse adaptive filtering strategies

- ▶ For each tested methods, tuning parameters optimized manually.
- ▶ S3MG leads to the minimal estimation error, and benefits from good tracking properties .

# Conclusion

## Solution to online penalized MSE problems

⇒ Proposition of a novel stochastic MM subspace algorithm.

- ✓ No need for tuning up a stepsize (in particular, no decreasing condition on a stepsize)
- ✓ Can be used in an adaptive context thanks to a forgetting factor
- ✓ Proven convergence guarantees
- ✓ Analysis of the convergence rate and of the complexity
- ✓ Good numerical performance w.r.t. state-of-the art methods.

# Some references



E. Chouzenoux and J.-C. Pesquet.

*A Stochastic Majorize-Minimize Subspace Algorithm for Online Penalized Least Squares Estimation.*  
Technical report, 2015. Available at <http://arxiv.org/abs/1512.08722>.



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.*  
IEEE Journal of Selected Topics in Signal Processing, Vol. 10, No. 2, pp. 224-241, March 2016.



E. Chouzenoux, A. Jezierska, J.-C. Pesquet and H. Talbot.

*A Majorize-Minimize Subspace Approach for  $\ell_2 - \ell_0$  Image Regularization.*  
SIAM Journal on Imaging Science, Vol. 6, No. 1, pp. 563-591, 2013.



E. Chouzenoux, J. Idier and S. Moussaoui.

*A Majorize-Minimize Strategy for Subspace Optimization Applied to Image Restoration.*  
IEEE Transactions on Image Processing, Vol. 20, No. 18, pp. 1517-1528, June 2011.

## Some references



E. Chouzenoux and J.-C. Pesquet.

*A Stochastic Majorize-Minimize Subspace Algorithm for Online Penalized Least Squares Estimation.*  
Technical report, 2015. Available at <http://arxiv.org/abs/1512.08722>.



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.*  
IEEE Journal of Selected Topics in Signal Processing, Vol. 10, No. 2, pp. 224-241, March 2016.



E. Chouzenoux, A. Jezierska, J.-C. Pesquet and H. Talbot.

*A Majorize-Minimize Subspace Approach for  $\ell_2 - \ell_0$  Image Regularization.*  
SIAM Journal on Imaging Science, Vol. 6, No. 1, pp. 563-591, 2013.



E. Chouzenoux, J. Idier and S. Moussaoui.

*A Majorize-Minimize Strategy for Subspace Optimization Applied to Image Restoration.*  
IEEE Transactions on Image Processing, Vol. 20, No. 18, pp. 1517-1528, June 2011.

Thank you !