

Fast Incremental Expectation Maximization: how many iterations for an ϵ -stationary point ?

Gersende Fort

CNRS & Institut de Mathématiques de Toulouse, France



CIRM "Optimization for Machine Learning", March 2020.

Based on a joint work with

- Pierre Gach - (IMT, Univ. Paul Sabatier, France)
- Eric Moulines - (CMAP, Ecole Polytechnique, France)

Acknowledgments:

- *Fondation Simone et Cino Del Duca* - **O**ptimisation et **S**imulation **M**onte Carlo : **E**ntrelacements
- *Labex CIMI* - **S**tochastic **D**escent Algorithms **w**ith **M**arkovian **I**nputs



Optimization problem for Statistical Learning

*An optimization problem occurring, for example,
in large scale Statistical inference*

The optimization problem

$$\operatorname{argmin}_{\theta \in \Theta} F(\theta), \quad F(\theta) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\theta) + R(\theta)$$

with

- $\Theta \subseteq \mathbb{R}^d$

- n is large

- $\mathcal{L}_i : \Theta \rightarrow \mathbb{R}$ is not explicit and of the form

$$\mathcal{L}_i(\theta) \stackrel{\text{def}}{=} -\log \int_{\mathcal{Z}} \exp(\langle s_i(z), \phi(\theta) \rangle) \, d\mu(z)$$

- **No convexity** assumptions

- (for the cvg analysis) Regularity properties of \mathcal{L}_i, R

Motivation: Statistical inference (1/3)

$$\operatorname{argmin}_{\theta \in \Theta} F(\theta), \quad F(\theta) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\theta) + R(\theta), \quad \mathcal{L}_i(\theta) \stackrel{\text{def}}{=} -\log \int_{\mathcal{Z}} \exp(\langle s_i(z), \phi(\theta) \rangle) d\mu(z)$$

- n observations Y_1, \dots, Y_n : a parametric statistical model indexed by $\theta \in \Theta$
- $R(\theta)$: a penalty term, a regularization term, a prior (in a Bayesian setting solved by MAP)
- A loss function associated to each Y_i : $\mathcal{L}_i(\theta)$
- $\exp(-\mathcal{L}_i)$ is a "sum" over **latent** variables $z \in \mathcal{Z}$. For example,

$$\mathcal{L}_i(\theta) = \underbrace{-}_{\text{negative}} \underbrace{\log}_{\text{log-}} \underbrace{\int_{\mathcal{Z}} \dots d\mu}_{\text{likelihood of } Y_i}$$

Example: mixture models (2/3)

- data: y_1, \dots, y_n modeled as iid from: $\sum_{\ell=1}^L \omega_{\ell} \mathcal{N}(\mu_{\ell}, 1)$ $\theta = (\omega_{1:L}, \mu_{1:L})$
- Equivalently: $z_i \in \{1, \dots, L\}$ s.t. $\mathcal{L}(Y_i | Z_i = \ell) \sim \mathcal{N}(\mu_{\ell}, 1)$ $\mathcal{L}(Z_i) \sim (\omega_{\ell})_{\ell}$.
- Joint distribution of (Y_i, Z_i) $\omega_{\ell}/\omega_L = \exp(\alpha_{\ell})$

$$\begin{aligned} \sum_{\ell=1}^L \mathbf{1}_{z_i=\ell} \omega_{\ell} \exp\left(-\frac{(y_i - \mu_{\ell})^2}{2}\right) &= \sum_{\ell=1}^L \mathbf{1}_{z_i=\ell} \exp\left(\ln \omega_{\ell} - \frac{(y_i - \mu_{\ell})^2}{2}\right) \\ &= \exp\left(\sum_{\ell=1}^L \mathbf{1}_{z_i=\ell} \left\{ \ln \omega_{\ell} - \frac{(y_i - \mu_{\ell})^2}{2} \right\}\right) \\ &= \exp\left(\sum_{\ell=1}^{L-1} \mathbf{1}_{z_i=\ell} \left\{ \alpha_{\ell} - \frac{(y_i - \mu_{\ell})^2}{2} + \frac{(y_i - \mu_L)^2}{2} \right\} - \ln\left(1 + \sum_{\ell=1}^{L-1} \exp(\alpha_{\ell})\right) - \frac{(y_i - \mu_L)^2}{2}\right) \end{aligned}$$

$$s_i(z) \stackrel{\text{def}}{=} (\mathbf{1}_{z=1}, \dots, \mathbf{1}_{z=L-1}, y_i \mathbf{1}_{z=1}, \dots, y_i \mathbf{1}_{z=L-1}, y_i, -1),$$

$$\phi(\theta) \stackrel{\text{def}}{=} \left(\alpha_1 - \frac{\mu_1^2 - \mu_L^2}{2}, \dots, \alpha_{L-1} - \frac{\mu_{L-1}^2 - \mu_L^2}{2}, \mu_1, \dots, \mu_{L-1}, \mu_L, \ln\left(1 + \sum_{\ell=1}^{L-1} \exp(\alpha_{\ell})\right) + \frac{\mu_L^2}{2} \right)$$

Example: Logistic regression (3/3)

- data: $y_1, \dots, y_n \in \{0, 1\}^n$ modeled as indep from

- $\mathcal{L}(Y_i|Z_i) \sim \text{Bern}\left(\frac{1}{1+\exp(-\eta_i(z_i))}\right)$ $\mathcal{L}(Z_i) \sim \mathcal{N}_p(\theta, \mathbf{I})$.

- The joint distribution of (Y_i, Z_i)

$$\left(\frac{\exp(\eta_i(z_i))}{1 + \exp(\eta_i(z_i))}\right)^{Y_i} \left(\frac{1}{1 + \exp(\eta_i(z_i))}\right)^{1-Y_i} \exp(-\|z_i - \theta\|^2/2) = \frac{\exp(Y_i \eta_i(z_i))}{1 + \exp(\eta_i(z_i))} \exp(-\|z_i - \theta\|^2/2).$$

- The likelihood of Y_i

$$\int_{\mathbb{R}^p} \frac{\exp(Y_i \eta_i(z))}{1 + \exp(\eta_i(z))} \exp(-\|z_i - \theta\|^2/2) dz = \int_{\mathcal{Z}} \exp(\langle s_i(z), \phi(\theta) \rangle) dz$$

$$s_i(z) = \left(y_i \eta_i(z) - \ln(1 + \exp(\eta_i(z))) - \|z\|^2/2, 1, z' \right)$$

$$\phi(\theta) = \left(1, -\|\theta\|^2/2, \theta' \right)$$

Which numerical tool ?

$$\operatorname{argmin}_{\theta \in \Theta} F(\theta), \quad F(\theta) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\theta) + R(\theta),$$
$$\mathcal{L}_i(\theta) \stackrel{\text{def}}{=} -\log \int_{\mathcal{Z}} \exp(\langle s_i(z), \phi(\theta) \rangle) \mu(\mathrm{d}z).$$

An algorithmic solution designed for

- large n : rare computations of a sum over n terms allowed,
- non convex setting

Solution: Expectation Maximization-based methods: Dempster et al. (1977), Wu (1983)

II- Expectation Maximization (EM) algorithms

EM algorithm: its derivation for this optim pbm,
its intractability,
an alternative.

EM: A Majorize-Minimization algorithm (1/2)

$$\text{Argmin}_{\theta} F(\theta) \quad F(\theta) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\theta) + R(\theta) \quad \mathcal{L}_i(\theta) = -\log \int_{\mathcal{Z}} \exp(\langle s_i(z), \phi(\theta) \rangle) d\mu(z).$$

- The surrogate function at the current point $\theta_k \in \Theta$, (Jensen's inequality)

$$\frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\theta) \leq \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\theta_k) + \langle \bar{s}(\theta_k), \phi(\theta_k) \rangle - \langle \bar{s}(\theta_k), \phi(\theta) \rangle$$

where

$$\bar{s}(\theta_k) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \bar{s}_i(\theta_k), \quad \bar{s}_i(\theta_k) \stackrel{\text{def}}{=} \int_{\mathcal{Z}} s_i(z) \frac{\exp(\langle s_i(z), \phi(\theta_k) \rangle)}{\exp(-\mathcal{L}_i(\theta_k))} \mu(dz).$$

- [E-step] Compute $\bar{s}(\theta_k)$
- [M-step] Minimize the majorizing function (under hyp: unique argmin)

$$\theta_{k+1} = T \circ \bar{s}(\theta_k) \stackrel{\text{def}}{=} \text{Argmin}_{\theta \in \Theta} \{-\langle \bar{s}(\theta_k), \phi(\theta) \rangle + R(\theta)\}.$$

EM in the Statistic-space (2/2)

θ_k	<div style="text-align: center; margin-bottom: 5px;">E</div> $\bar{s}(\theta_k)$	<div style="text-align: center; margin-bottom: 5px;">M</div> $\theta_{k+1} = T \circ \bar{s}(\theta_k)$	<div style="text-align: center; margin-bottom: 5px;">E</div> $\bar{s} \circ T \circ \bar{s}(\theta_k)$	<div style="text-align: center; margin-bottom: 5px;">M</div> $\theta_{k+2} = T \circ \bar{s} \circ T \circ \bar{s}(\theta_k)$	$F(\theta)$
s_k	s_k	$T(s_k)$	$s_{k+1} = \bar{s} \circ T(s_k)$	$T \circ \bar{s} \circ T(s_k)$	$V(s)$

We will see EM-based algorithms as evolving in the "*s*-space"; the objective function: $V = F \circ T$

- If convergence, to the roots of

$$h(s) \stackrel{\text{def}}{=} \bar{s} \circ T(s) - s.$$

Under assumptions, the roots of h are the roots of \dot{V} . Do Stochastic EM avoid traps ?

not the topic today

EM is designed to find the roots of h

From EM to Stochastic Approximation within EM

- Each iteration of EM requires

$$s_{k+1} = \frac{1}{n} \sum_{i=1}^n \bar{s}_i(T(s_k))$$

It does not answer the specifications for large scale learning.

- It is designed to find the zeros of h and h is intractable:

$$h(s) = \frac{1}{n} \sum_{i=1}^n \bar{s}_i \circ T(s) - s = \mathbb{E} [\bar{s}_I \circ T(s) - s] \quad I \sim \mathcal{U}(\{1, \dots, n\})$$

what about *Stochastic Approximation within EM approaches* ?

among the many "Stochastic EM" algorithms [SEM] Celeux & Diebolt (1985); [MCEM] Wei & Tanner (1990), Fort & Moulines (2003); [SAEM] Delyon, Lavielle & Moulines (1999); Kuhn & Lavielle (2004); [Online EM] Cappé & Moulines (2009), [Incremental EM] Neal & Hinton (1999) ...

Stochastic Approximation (SA) within EM:

EM \rightarrow designed to find the roots of $h(s) \stackrel{\text{def}}{=} \bar{s} \circ T(s) - s$

- What is SA ?

Solve $h(s) = 0$ when $h(s) \stackrel{\text{def}}{=} \mathbb{E}[H(U, s)]$ by:

$$\hat{S}^{k+1} = \hat{S}^k + \gamma_{k+1} \underbrace{H(U_{k+1}, \hat{S}^k)}_{\text{approx. of } h(\hat{S}^k)}, \quad U_{k+1} \sim U \text{ or "almost"}$$

- Here, possibly **two** sources of intractability

$$h(s) = \frac{1}{n} \sum_{i=1}^n \bar{s}_i \circ T(s) - s = \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{Z}} s_i(z) \pi(\mathrm{d}z; T(s)) - s$$

Delyon et al. (1999), the second intractability only with i.i.d. U'_k .

Extensions to MCMC sampling, first by Lavielle & Kuhn (2004)

1st idea: SA

$$h(s) = \frac{1}{n} \sum_{i=1}^n \bar{s}_i \circ T(s) - s = \mathbb{E} [\bar{s}_J \circ T(s) - s] \quad J \sim \mathcal{U}([n])$$

Instead of the EM iterations $s_{k+1} = h(s_k) = \bar{s} \circ T(s_k)$ do:

Data: $\hat{S}^0 \in \mathcal{S}$, $\gamma_k \in (0, \infty)$ for $k \geq 1$

Result: The SA sequence: $\hat{S}^k, k = 0, \dots,$

for $k \geq 1$ **do**

 Sample J_{k+1} uniformly on $[n]$;

$$\hat{S}^{k+1} = \hat{S}^k + \gamma_{k+1} (\bar{s}_{J_{k+1}} \circ T(\hat{S}^k) - \hat{S}^k)$$

Fast Incremental EM (FIEM) [FIEM] Karimi et al. (2019); [SAGA] Defazio et al. (2014)

$$h(s) = \mathbb{E} \left[\bar{s}_I \circ T(s) - s + V_{k+1} \right] \quad I \sim \mathcal{U}([n]), \quad \mathbb{E}[V_{k+1}] = 0$$

Data: $\hat{S}^0 \in \mathcal{S}$, $\gamma_k \in (0, \infty)$ for $k \geq 1$

Result: The FIEM sequence: $\hat{S}^k, k = 0, \dots,$

$S_{0,i} = \bar{s}_i \circ T(\hat{S}^0)$ for all $i \in [n]$;

$$\tilde{S}^0 = n^{-1} \sum_{i=1}^n S_{0,i};$$

for $k \geq 1$ **do**

* Sample I_{k+1} uniformly on $[n]$;

* $S_{k+1,i} = S_{k,i}$ for $i \neq I_{k+1}$;

* $S_{k+1,I_{k+1}} = \bar{s}_{I_{k+1}} \circ T(\hat{S}^k)$;

* $\tilde{S}^{k+1} = \tilde{S}^k + n^{-1} (S_{k+1,I_{k+1}} - S_{k,I_{k+1}})$;

Sample J_{k+1} uniformly on $[n]$;

$\hat{S}^{k+1} = \hat{S}^k + \gamma_{k+1} (\bar{s}_{J_{k+1}} \circ T(\hat{S}^k) - \hat{S}^k - \{S_{k+1,J_{k+1}} - \tilde{S}^{k+1}\})$

* iterative comput. of the sum:

$$\tilde{S}^{k+1} = \frac{1}{n} \sum_{i=1}^n \bar{s}_i \circ T(\hat{S}^{\text{last time when } \#i \text{ sampled}})$$

III- Non asymptotic convergence bounds

Non asymptotic \rightarrow a maximal number of iterations

How to define the estimate of the solution ?

“Convergence bounds”: in which sense ?

How to choose the stepsize sequence ?

Non asymptotic convergence bounds: which ones ?

- The user chooses a maximal length of the path: K_{\max} .

K is a random stopping time in the range $\{0, \dots, K_{\max} - 1\}$. Ghadimi & Lan (2013)

- mean "distance" to the roots of the gradient of $V \stackrel{\text{def}}{=} F \circ T$

$$\inf_{k \leq K_{\max}} \mathbb{E} \left[\|\dot{V}(\hat{S}^k)\|^2 \right] \leq \mathbb{E}_g \stackrel{\text{def}}{=} \mathbb{E} \left[\|\dot{V}(\hat{S}^K)\|^2 \right]$$

- mean "distance" to the roots of h

$$\mathbb{E}_h \stackrel{\text{def}}{=} \mathbb{E} \left[\|h(\hat{S}^K)\|^2 \right]$$

- Under the stated assumptions

$$\frac{1}{v_{\max}^2} \mathbb{E}_g \leq \mathbb{E}_h$$

Assumptions

A1 $\Theta \subseteq \mathbb{R}^d$ is an open set. (Z, \mathcal{Z}) is a measurable space and μ is a σ -finite positive measure on Z . The functions $\phi : \Theta \rightarrow \mathbb{R}^q$, $s_i : Z \rightarrow \mathbb{R}^q$ for all $i \in [[n]]$ and $R : \Theta \rightarrow \mathbb{R}$ are measurable functions. Finally, for any $\theta \in \Theta$ and $i \in [[n]]$, $-\infty < \mathcal{L}_i(\theta) < \infty$.

A2 For all $\theta \in \Theta$ and $i \in [[n]]$, the expectation $\bar{s}_i(\theta)$ exists. For any $s \in \mathbb{R}^q$, $\text{Argmin}_{\theta \in \Theta} (-\langle s, \phi(\theta) \rangle + R(\theta))$ is a (non empty) singleton denoted by $\{T(s)\}$.

A3 ϕ, R are C^1 on Θ . T is C^1 on \mathbb{R}^q .

For any $s \in \mathbb{R}^q$, $B(s) \stackrel{\text{def}}{=} \nabla(\phi \circ T)(s)$ is a symmetric $q \times q$ matrix and there exist $0 < v_{\min} \leq v_{\max} < \infty$ such that for all $s \in \mathcal{S}$, the spectrum of $B(s)$ is in $[v_{\min}, v_{\max}]$.

For any $i \in [[n]]$, $\bar{s}_i \circ T$ is globally Lipschitz on \mathbb{R}^q with constant L_i .

$s \mapsto B^T(s) (\bar{s} \circ T(s) - s)$ is globally Lipschitz on \mathbb{R}^q with constant $L_{\dot{V}}$.

Corollary: R is C^1 ; \mathcal{L}_i is C^1 ; \dot{V} is Lipschitz; $\|\dot{V}(s)\| \leq v_{\max} \|h(s)\|$.

Sketch of proof, $V \stackrel{\text{def}}{=} F \circ T$

An upper bound of $\mathbb{E} [\|h(\hat{S}^K)\|^2]$ when $\hat{S}^{k+1} = \hat{S}^k + \gamma_{k+1} H_{k+1}$

- A Taylor expansion, first order + Gradient Lipschitz

$$V(\hat{S}^{k+1}) \leq V(\hat{S}^k) + \gamma_{k+1} \langle H_{k+1}, \dot{V}(\hat{S}^k) \rangle + \gamma_{k+1}^2 \frac{L_{\dot{V}}}{2} \|H_{k+1}\|^2$$

- The expectation

$$\mathbb{E} [V(\hat{S}^{k+1})] \leq \mathbb{E} [V(\hat{S}^k)] + \gamma_{k+1} \mathbb{E} [\langle h(\hat{S}^k), \dot{V}(\hat{S}^k) \rangle] + \gamma_{k+1}^2 \frac{L_{\dot{V}}}{2} \mathbb{E} [\|H_{k+1}\|^2]$$

- The assumption (Lyapunov contraction)

$$\mathbb{E} [V(\hat{S}^{k+1})] \leq \mathbb{E} [V(\hat{S}^k)] - \gamma_{k+1} v_{\min} \mathbb{E} [\|h(\hat{S}^k)\|^2] + \gamma_{k+1}^2 \frac{L_{\dot{V}}}{2} \mathbb{E} [\|H_{k+1}\|^2]$$

- A sum from $k = 0$ to $k = K_{\max} - 1$,

$$v_{\min} \sum_{k=0}^{K_{\max}-1} \gamma_{k+1} \mathbb{E} [\|h(\hat{S}^k)\|^2] \leq \mathbb{E} [V(\hat{S}^0)] - \mathbb{E} [V(\hat{S}^{K_{\max}})] + \frac{L_{\dot{V}}}{2} \sum_{k=0}^{K_{\max}-1} \gamma_{k+1}^2 \mathbb{E} [\|H_{k+1}\|^2]$$

- Few pages later:

$$\sum_{k=0}^{K_{\max}-1} \alpha_k \mathbb{E} [\|h(\hat{S}^k)\|^2] \leq \mathbb{E} [V(\hat{S}^0)] - \mathbb{E} [V(\hat{S}^{K_{\max}})]$$

Result 1 (Gach, F., Moulines-2020)

Assume A1 to A3 and set $L^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n L_i^2$. Choose $\mu \in (0, 1)$.

Let K be a $\{0, \dots, K_{\max} - 1\}$ -valued uniform r.v.

Run FIEM with a **constant** step size

$$\gamma_\ell = \frac{\sqrt{C}}{n^{2/3}L}$$

where $C \in (0, 1)$ is the unique solution of "an equation depending upon $v_{\min}, L, L_{\dot{Y}}, \mu$ "

Then, for any $n \geq 2$ and $K_{\max} \geq 1$

$$E_h \leq \frac{n^{2/3}}{K_{\max}} \frac{L}{\sqrt{C}(1-\mu)v_{\min}} \mathbb{E} \left[V(\hat{S}^0) - V(\hat{S}^{K_{\max}}) \right].$$

Corollaries of Result 1

$$\mathbb{E}_h \leq \frac{n^{2/3}}{K_{\max}} \frac{L}{\sqrt{C}(1-\mu)v_{\min}} \mathbb{E} \left[V(\hat{S}^0) - V(\hat{S}^{K_{\max}}) \right] \quad \gamma_k = \frac{\sqrt{C}}{n^{2/3}L}$$

- Dependence upon n and K_{\max} : as in Karimi et al. (2019)
- Constant stepsize $O(n^{-2/3})$ as in Karimi et al. (2019)
- Precision ε : $K_{\max} = M n^{2/3} \varepsilon^{-1}$
- $C \in (0, 1)$ is explicit \rightarrow "definition" of $\gamma_k \rightarrow$ improves on previous results

$$\sqrt{C} \frac{L\check{v}}{2L} \left(\frac{1}{n^{2/3}} + \frac{C}{1/2 - Cn^{-1/3}} \left(\frac{1}{n} + 2 \right) \right) = \mu v_{\min}$$

when $n \rightarrow \infty$ $C^{3/2} = \frac{L}{2L\check{v}} \mu v_{\min}$

Result 2 (Gach, F., Moulines-2020)

Assume A1 to A3 and set $L^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n L_i^2$. Choose $\mu \in (0, 1)$.

Let K be a $\{0, \dots, K_{\max} - 1\}$ -valued uniform r.v.

Run FIEM with a **constant** step size

$$\gamma_\ell = \frac{\sqrt{C}}{n^{1/3} K_{\max}^{1/3} L}$$

where $C > 0$ is the unique solution of "an equation depending upon $v_{\min}, L, L_{\check{Y}}, \mu$ "

Then, for any $n \geq 1$ and $K_{\max} \geq 1$

$$E_h \leq \frac{n^{1/3}}{K_{\max}^{2/3}} \frac{L}{\sqrt{C}(1-\mu)v_{\min}} \mathbb{E} \left[V(\hat{S}^0) - V(\hat{S}^{K_{\max}}) \right].$$

Corollaries of Result 2

$$E_h \leq \frac{n^{1/3}}{K_{\max}^{2/3}} \frac{L}{\sqrt{C}(1-\mu)v_{\min}} \mathbb{E} \left[V(\hat{S}^0) - V(\hat{S}^{K_{\max}}) \right] \quad \gamma_\ell = \frac{\sqrt{C}}{n^{1/3}K_{\max}^{1/3}L}.$$

- Dependence upon n and K_{\max} : new result.
- Constant stepsize $O(n^{-1/3}K_{\max}^{-1/3})$
- Precision ε : $K_{\max} = M \sqrt{n} \varepsilon^{-3/2} \rightarrow \text{stepsize} = O(\sqrt{\varepsilon}/\sqrt{n})$
- $C \in (0, 1)$ is explicit \rightarrow "definition" of γ_k

Result 3 (Gach, F., Moulines-2020)

Let K be a $\{0, \dots, K_{\max} - 1\}$ -valued r.v. with weights $p_0, \dots, p_{K_{\max}-1}$, $\inf_k p_k > 0$.

Let $C \in (0, 1)$ be the unique solution of "an equation depending upon $v_{\min}, L, L_{\dot{V}}, \mu$ " Run FIEM with

$$\gamma_{k+1} \stackrel{\text{def}}{=} \frac{g_k}{n^{2/3}L}, \quad g_k \stackrel{\text{def}}{=} F_{n,C}^{-1} \left(\frac{p_k}{\max_{\ell} p_{\ell}} \frac{v_{\min} \sqrt{C}}{2L} \frac{1}{n^{2/3}} \right);$$

$$F_{n,C} : x \mapsto \frac{1}{Ln^{2/3}} x (v_{\min} - x f_n(C)), \quad f_n(C) \stackrel{\text{def}}{=} \frac{L_{\dot{V}}}{2L} \left(\frac{1}{n^{2/3}} + \frac{C}{1/2 - Cn^{-1/3}} \left(\frac{1}{n} + 2 \right) \right).$$

For any $n \geq 2$ and $K_{\max} \geq 1$, we have

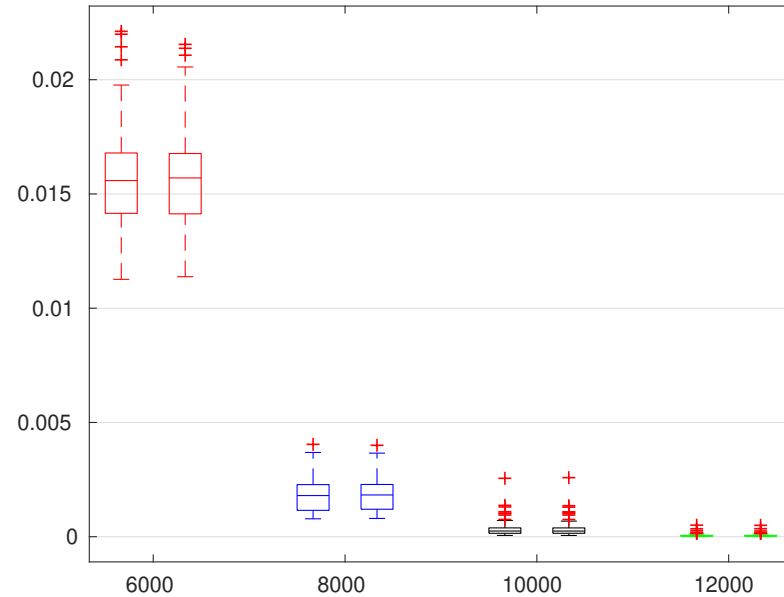
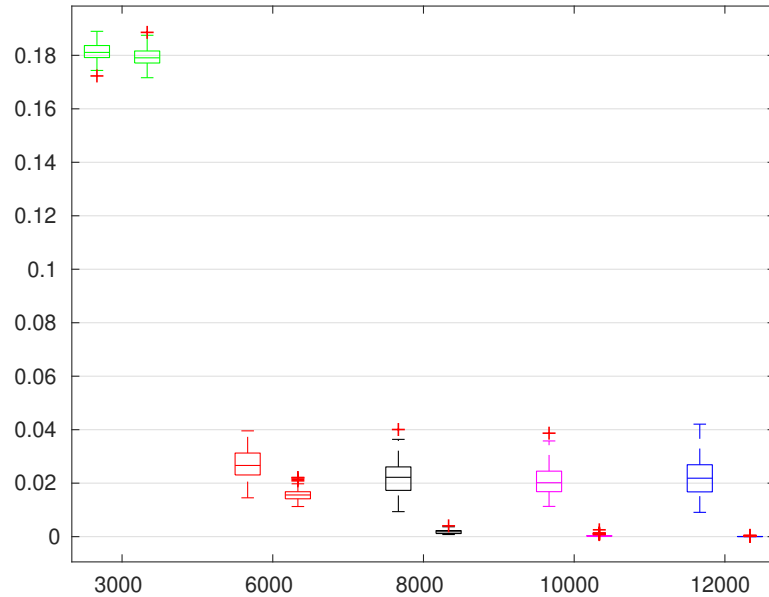
$$E_h \leq n^{2/3} \max_k p_k \frac{2L}{\sqrt{C}v_{\min}} \mathbb{E} \left[V(\hat{S}^0) - V(\hat{S}^{K_{\max}}) \right].$$

-
- Optimal sampling strategy : $p_k = 1/K_{\max}$ \rightarrow same as Result 1.
 - New result.

On a toy example (1/4)

- Toy example: $\text{Argmin}_{\theta} F(\theta) = \theta_{\star}$ is explicit.
- Compare SA and FIEM (left) through boxplots of $\|\theta_k - \theta_{\star}\|$ for different k ; same thing for FIEM and FIEM-opt in the convergence phase (right)

$$\text{FIEM-opt } \hat{S}^{k+1} = \hat{S}^k + \gamma_{k+1} (T_{k+1} + \lambda_{k+1} V_{k+1})$$

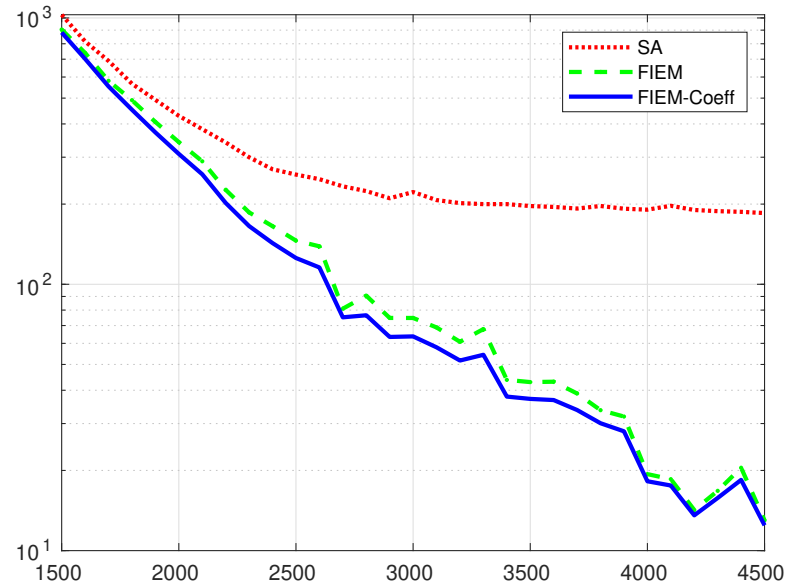
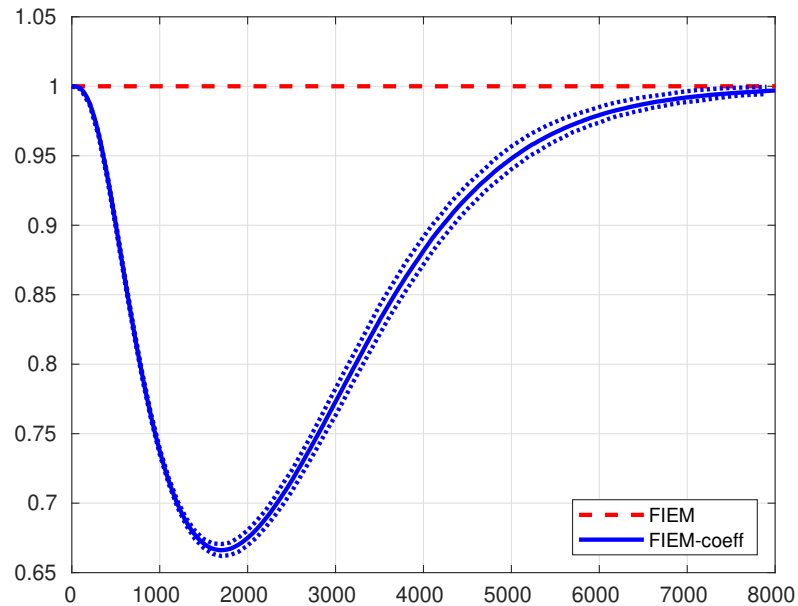


On a toy example (2/4)

- FIEM and FIEM-opt

(left) The coefficient $k \mapsto \lambda_{k+1}$;

(right) The L^2 -moment of the field estimated by Monte Carlo: $\mathbb{E} [\|H_{k+1}\|^2]$.

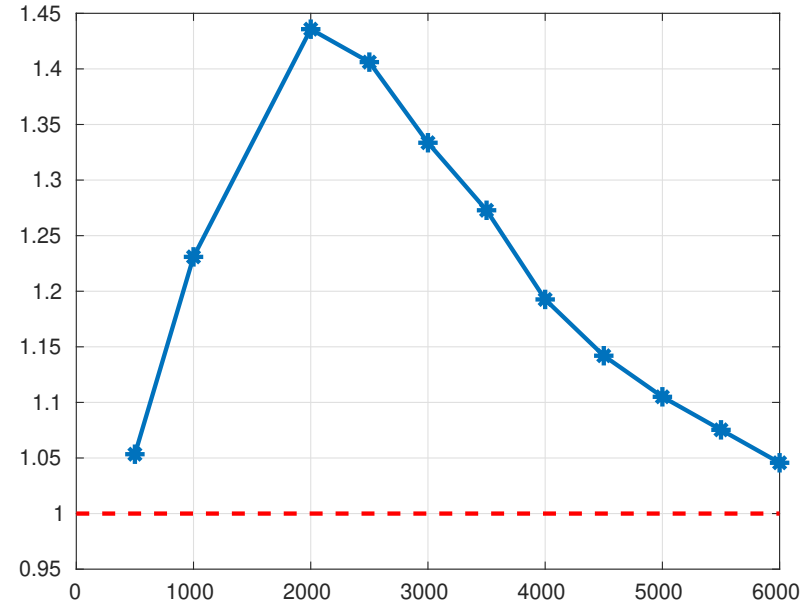
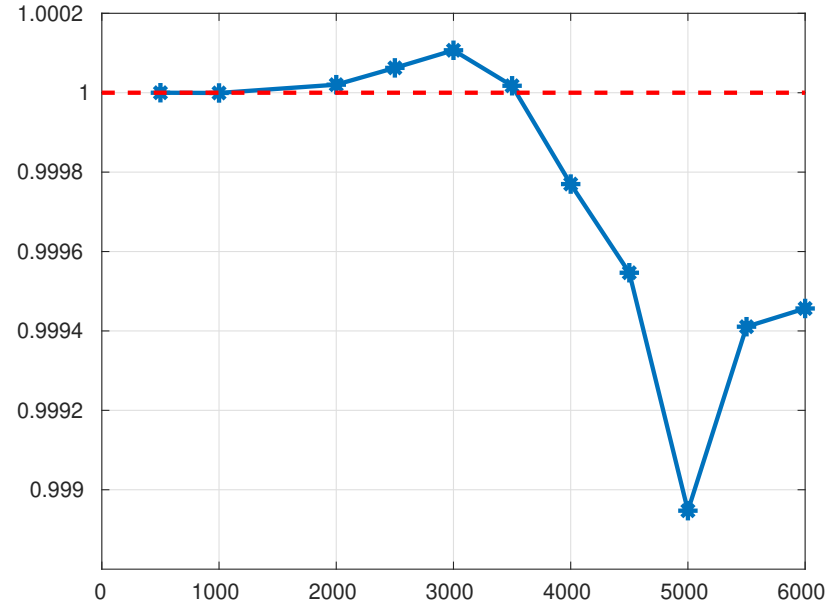


On a toy example (3/4)

- FIEM and FIEM-opt - Monte Carlo evaluation of $\mathbb{E} [\|\theta_k - \theta_\star\|]$ and $\text{std}(\|\theta_k - \theta_\star\|)$.

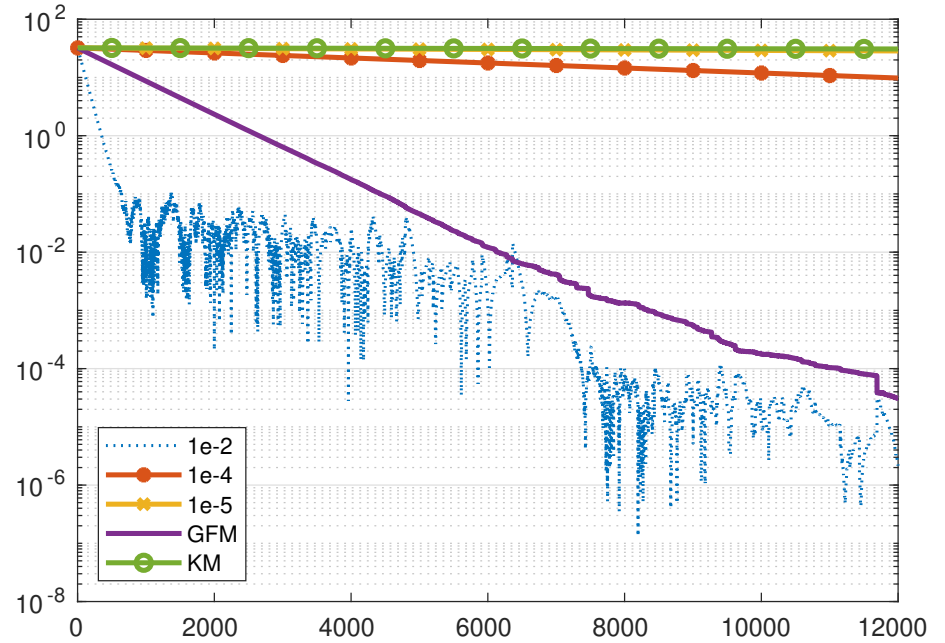
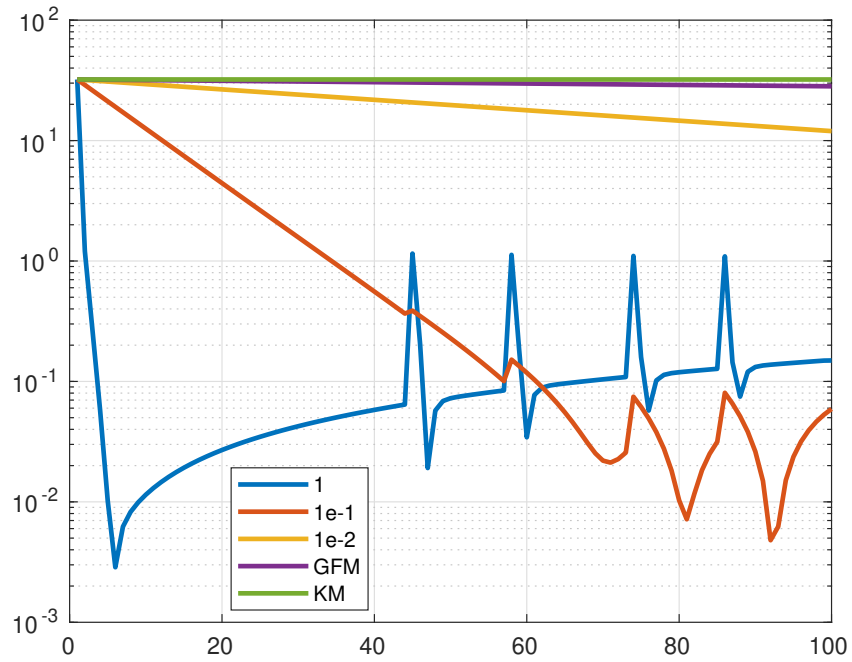
(left) ratio of the expectations FIEM / FIEM-opt

(right) ratio of the std FIEM / FIEM-opt



On a toy example (4/4)

- For FIEM, constant stepsize γ , observe $k \mapsto \|\theta_k - \theta_\star\|$



IV- What happens when \bar{s}_i is not explicit ?

The perturbed case: when the expectation are approximated, which conditions on the approximations in order to have the same rates as in the "exact" case?

The Perturbed FIEM algorithm

Data: $\hat{S}^0 \in \mathcal{S}$, $\gamma_k \in (0, \infty)$ for $k \geq 1$

Result: The FIEM sequence: $\hat{S}^k, k = 0, \dots,$

$S_{0,i} = \tilde{s}_i$, an approximation of $\bar{s}_i \circ T(\hat{S}^0)$ for all $i \in [[n]]$;

$$\tilde{S}^0 = n^{-1} \sum_{i=1}^n S_{0,i};$$

for $k \geq 1$ **do**

* Sample I_{k+1} uniformly on $[[n]]$;

* $S_{k+1,i} = S_{k,i}$ for $i \neq I_{k+1}$;

* $S_{k+1,I_{k+1}} = \check{S}_{k+1}$ an approximation of $\bar{s}_{I_{k+1}} \circ T(\hat{S}^k)$;

* $\tilde{S}^{k+1} = \tilde{S}^k + n^{-1} (S_{k+1,I_{k+1}} - S_{k,I_{k+1}})$;

Sample J_{k+1} uniformly on $[[n]]$;

Compute \tilde{s}_{k+1} , an approximation of $\bar{s}_{J_{k+1}} \circ T(\hat{S}^k)$;

$$\hat{S}^{k+1} = \hat{S}^k - \gamma_{k+1} (\tilde{s}_{k+1} - \hat{S}^k - \{S_{k+1,J_{k+1}} - \tilde{S}^{k+1}\})$$

Assumptions (on the perturbations)

There exist positive sequences $\{m_k, k \geq 0\}$ and $\{\bar{m}_k, k \geq 0\}$, positive numbers $M^{(1)}$ and $M^{(2)}$ and $M_\nu^{(2)} \geq 0$ such that for all $k \geq 0$, the approximations \tilde{s}_{k+1} and \check{S}_{k+1} satisfy

$$\mathbb{E}[\|\check{S}_{k+1} - \bar{s}_{I_{k+1}} \circ T(\hat{S}^k)\|^2] \leq \frac{M^{(1)}}{\bar{m}_{k+1}},$$

$$\mathbb{E}[\|\mathbb{E}[\tilde{s}_{k+1} - \bar{s}_{J_{k+1}} \circ T(\hat{S}^k) | J_{k+1}, I_{k+1}, \hat{S}^k]\|^2] \leq \frac{M_\nu^{(2)}}{m_{k+1}^2},$$

$$\mathbb{E}[\|\tilde{s}_{k+1} - \bar{s}_{J_{k+1}} \circ T(\hat{S}^k)\|^2] \leq \frac{M^{(2)}}{m_{k+1}}.$$

Well ... in a Monte Carlo approximation

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|\frac{1}{\bar{m}_{k+1}} \sum_{l=1}^{\bar{m}_{k+1}} s_i(Z_{l,k}) - \bar{s}_i \circ T(\hat{S}^k)\|^2] \leq \frac{\mathbb{E}[C(\hat{S}^k)]}{\bar{m}_{k+1}} \leq \frac{M^{(1)}}{\bar{m}_{k+1}},$$

Convergence result (Gach, F., Moulines (2020))

$$\begin{aligned} E_h &\leq \frac{n^{2/3}}{K_{\max}} C_0 \\ &+ \frac{C_1}{n^{2/3}} \left\{ 1 \wedge \frac{n}{K_{\max}} \right\} \mathbb{E} \left[\varepsilon^{(0)} \right] \quad \text{error when initializing} \\ &+ \frac{C_1}{n^{5/3}} \left\{ 1 \wedge \frac{n}{K_{\max}} \right\} \sum_{k=0}^{K_{\max}-1} \mathbb{E} \left[\|\eta_{k+1}^{(1)}\|^2 \right] \quad \text{Error: on the control variate} \\ &+ C_2 \frac{1}{K_{\max}} \sum_{k=0}^{K_{\max}-1} \mathbb{E} \left[\|\mathbb{E} \left[\eta_{k+1}^{(2)} \mid \mathcal{F}_k \right]\|^2 \right] \quad \eta^{(2)}: \text{second error} \\ &+ \frac{C_1}{2(1+\nu)} \frac{1}{n^{2/3} K_{\max}} \sum_{k=0}^{K_{\max}-1} \left(\mathbb{E} \left[\|\eta_{k+1}^{(2)}\|^2 \right] + \mathbb{E} \left[\|\mathbb{E} \left[\eta_{k+1}^{(2)} \mid \mathcal{F}_{k+3/4} \right]\|^2 \right] \right); \end{aligned}$$

In the case of a Monte Carlo approximation

For a precision ε ,

- with an **unbiased** Monte Carlo approximation (ex. i.i.d.)

$$+ K_{\max} = Mn^{2/3} \varepsilon^{-1}$$

$$+ m = n^{-2/3} \varepsilon^{-1} \quad \bar{m} = (n^{-2/3} \varepsilon^{-1}) \wedge (n^{-1} \varepsilon^{-2})$$

- with a **biased** Monte Carlo approximation (ex. Markov chain Monte Carlo)

$$+ K_{\max} = Mn^{2/3} \varepsilon^{-1}$$

$$+ m = (n^{-2/3} \varepsilon^{-1}) \vee \varepsilon^{-1/2} \quad \bar{m} = (n^{-2/3} \varepsilon^{-1}) \wedge (n^{-1} \varepsilon^{-2})$$