

Stochastic Levenberg-Marquardt Methods for Noisy Derivative-Free Optimization with Complexity Results and Application to Data Assimilation

Clément Royer (*Université Paris Dauphine-PSL*)

Joint work with Elhoucine Bergou, Youssef Diouane, Vyacheslav Kungurtsev

SIAM CSE, March 4, 2021



Then

- Stochastic Levenberg-Marquardt;
- Stochastic estimates+probabilistic models.

Now

- A new version, with more complexity results;
A stochastic Levenberg-Marquardt method using random models with complexity results and application to data assimilation.
- **Another work** on Levenberg-Marquardt methods with constraints (but derivatives...).
A nonmonotone matrix-free algorithm for nonlinear equality-constrained least-squares problems.

- 1 Derivative-free least squares
- 2 Introducing randomness
- 3 New complexity results

Motivation: data assimilation problem

Strongly constrained 4DVAR formulation

$$\min_{\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_T \in \mathbb{R}^n} \underbrace{\frac{1}{2} \|\mathbf{z}_0 - \mathbf{z}_b\|_{\mathbf{B}^{-1}}^2}_{\text{Background state}} + \underbrace{\frac{1}{2} \sum_{i=0}^T \|\mathbf{y}_i - \mathcal{H}(\mathbf{z}_i)\|_{\mathbf{R}_i^{-1}}^2}_{\text{Measurement fit}}$$

s.t. $\underbrace{\mathbf{z}_i = \mathcal{M}_i(\mathbf{z}_{i-1}), i = 1, \dots, T.}_{\text{Governing equations}}$

Motivation: data assimilation problem

Strongly constrained 4DVAR formulation

$$\min_{\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_T \in \mathbb{R}^n} \underbrace{\frac{1}{2} \|\mathbf{z}_0 - \mathbf{z}_b\|_{\mathbf{B}^{-1}}^2}_{\text{Background state}} + \underbrace{\frac{1}{2} \sum_{i=0}^T \|\mathbf{y}_i - \mathcal{H}(\mathbf{z}_i)\|_{\mathbf{R}_i^{-1}}^2}_{\text{Measurement fit}}$$

s.t.

$$\underbrace{\mathbf{z}_i = \mathcal{M}_i(\mathbf{z}_{i-1}), \quad i = 1, \dots, T.}_{\text{Governing equations}}$$

Reformulation ($\mathbf{x} \leftrightarrow \mathbf{z}_0$)

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) := \frac{1}{2} \left(\|\mathbf{x} - \mathbf{z}_b\|_{\mathbf{B}^{-1}}^2 + \|\mathbf{y} - \mathcal{H}(\mathbf{x})\|_{\mathbf{R}^{-1}}^2 \right).$$

- Derivatives expensive to compute;
- \mathbf{B} unknown \Rightarrow Estimated via random samples. instead!

A more general setup

Derivative-free nonlinear least-squares

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) := \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|^2$$

- $\mathbf{r} : \mathbb{R}^n \mapsto \mathbb{R}^m, \mathbf{r} \in \mathcal{C}^{1,1};$
- $\mathbf{J}(\mathbf{x}) := [\nabla r_i(\mathbf{x})^\top] \in \mathbb{R}^{m \times n}$ **not available**;
- Values of \mathbf{r} only accessed through **noisy estimates**.

A more general setup

Derivative-free nonlinear least-squares

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) := \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|^2$$

- $\mathbf{r} : \mathbb{R}^n \mapsto \mathbb{R}^m, \mathbf{r} \in \mathcal{C}^{1,1};$
- $\mathbf{J}(\mathbf{x}) := [\nabla r_i(\mathbf{x})^\top] \in \mathbb{R}^{m \times n}$ **not available**;
- Values of \mathbf{r} only accessed through **noisy estimates**.

Classical Levenberg-Marquardt approach

- Gauss-Newton model $f(\mathbf{x} + \mathbf{s}) \approx \frac{1}{2} \|\mathbf{r}(\mathbf{x}) + \mathbf{J}(\mathbf{x})\mathbf{s}\|^2 + \frac{\gamma}{2} \|\mathbf{s}\|^2;$
- **Regularization parameter** γ set adaptively.

Derivative-free Levenberg-Marquardt for $\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|^2$

Inputs: $\mathbf{x}_0 \in \mathbb{R}^n$, $\gamma_0 > 0$, $\eta > 0$.

Iteration j : Given (\mathbf{x}_j, γ_j) ,

- Compute $\mathbf{r}_{m_j} \approx \mathbf{r}(\mathbf{x}_j)$, $\mathbf{J}_{m_j} \approx \mathbf{J}(\mathbf{x}_j)$ and
$$\mathbf{s}_j \approx \operatorname{argmin}_{\mathbf{s}} m_j(\mathbf{s}) := \frac{1}{2} \|\mathbf{r}_{m_j} + \mathbf{J}_{m_j} \mathbf{s}\|^2 + \frac{\gamma_j \|\mathbf{J}_{m_j}^\top \mathbf{r}_{m_j}\|}{2} \|\mathbf{s}\|^2.$$
- Compute $\mathbf{r}_j^0 \approx \mathbf{r}(\mathbf{x}_j)$ and $\mathbf{r}_j^s \approx \mathbf{r}(\mathbf{x}_j + \mathbf{s}_j)$.
- If $\frac{\frac{1}{2} \|\mathbf{r}_j^0\|^2 - \frac{1}{2} \|\mathbf{r}_j^s\|^2}{m_j(0) - m_j(\mathbf{s})} \geq \eta$, set $\mathbf{x}_{j+1} = \mathbf{x}_j + \mathbf{s}_j$ and $\gamma_{j+1} = 0.5\gamma_j$;
- Otherwise, set $\mathbf{x}_{j+1} = \mathbf{x}_j$ and $\gamma_{j+1} = 2\gamma_j$.

Derivative-free Levenberg-Marquardt for $\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|^2$

Inputs: $\mathbf{x}_0 \in \mathbb{R}^n$, $\gamma_0 > 0$, $\eta > 0$.

Iteration j : Given (\mathbf{x}_j, γ_j) ,

- Compute $\mathbf{r}_{m_j} \approx \mathbf{r}(\mathbf{x}_j)$, $\mathbf{J}_{m_j} \approx \mathbf{J}(\mathbf{x}_j)$ and $\mathbf{s}_j \approx \operatorname{argmin}_{\mathbf{s}} m_j(\mathbf{s}) := \frac{1}{2} \|\mathbf{r}_{m_j} + \mathbf{J}_{m_j} \mathbf{s}\|^2 + \frac{\gamma_j \|\mathbf{J}_{m_j}^\top \mathbf{r}_{m_j}\|}{2} \|\mathbf{s}\|^2$.
- Compute $\mathbf{r}_j^0 \approx \mathbf{r}(\mathbf{x}_j)$ and $\mathbf{r}_j^s \approx \mathbf{r}(\mathbf{x}_j + \mathbf{s}_j)$.
- If $\frac{\frac{1}{2} \|\mathbf{r}_j^0\|^2 - \frac{1}{2} \|\mathbf{r}_j^s\|^2}{m_j(0) - m_j(\mathbf{s})} \geq \eta$, set $\mathbf{x}_{j+1} = \mathbf{x}_j + \mathbf{s}_j$ and $\gamma_{j+1} = 0.5\gamma_j$;
- Otherwise, set $\mathbf{x}_{j+1} = \mathbf{x}_j$ and $\gamma_{j+1} = 2\gamma_j$.

- Goal: **Complexity results**.
- Key: **Accuracy properties** for the models/estimates.

First-order accuracy

A model m_j is called $(\kappa_{ef}, \kappa_{eg})$ -accurate if

$$\|\mathbf{J}(\mathbf{x}_j)^\top \mathbf{r}(\mathbf{x}_j) - \mathbf{J}_{m_j}^\top \mathbf{r}_{m_j}\| \leq \frac{\kappa_{eg}}{\gamma_j} \quad \text{and} \quad \left| \frac{1}{2} \|\mathbf{r}(\mathbf{x}_j)\|^2 - \frac{1}{2} \|\mathbf{r}_{m_j}\|^2 \right| \leq \frac{\kappa_{ef}}{\gamma_j^2}.$$

Accurate function estimates

Estimates $(\mathbf{r}_j^0, \mathbf{r}_j^s)$ are ε_f -accurate if

$$\left| \frac{1}{2} \|\mathbf{r}_j^0\|^2 - \frac{1}{2} \|\mathbf{r}(\mathbf{x}_j)\|^2 \right| \leq \frac{\varepsilon_f}{\gamma_j^2} \quad \text{and} \quad \left| \frac{1}{2} \|\mathbf{r}_j^s\|^2 - \frac{1}{2} \|\mathbf{r}(\mathbf{x}_j + \mathbf{s}_j)\|^2 \right| \leq \frac{\varepsilon_f}{\gamma_j^2}$$

Complexity analysis without derivatives

Goal: Bound the number of iterations to reach an ϵ_d -point \mathbf{x} such that

$$\|\nabla f(\mathbf{x})\| = \|\mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x})\| \leq \epsilon_d$$

as a function of ϵ_d .

Complexity analysis without derivatives

Goal: Bound the number of iterations to reach an ϵ_d -point \mathbf{x} such that

$$\|\nabla f(\mathbf{x})\| = \|\mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x})\| \leq \epsilon_d$$

as a function of ϵ_d .

Theorem (BDKR '18)

Suppose that

- Every model m_j is $(\kappa_{ef}, \kappa_{eg})$ -accurate.
- Every estimate pair $(\mathbf{r}_j^0, \mathbf{r}_j^s)$ is ε_f -accurate.

Then, the method reaches an ϵ_d -point in at most

$$\mathcal{O}\left(\kappa^2 \epsilon_d^{-2}\right)$$

iterations with $\kappa = \max\{\kappa_{ef}, \kappa_{eg}, \varepsilon_f\}$.

- Derivative-based case: $\mathcal{O}(\epsilon_d^{-2})$.

Considering randomness

Recall : Our function estimates are noisy.

Inputs: $\mathbf{x}_0 \in \mathbb{R}^n$, $\gamma_0 > 0$.

Iteration j : Given (\mathbf{x}_j, γ_j) ,

- Compute $\mathbf{r}_{m_j} \approx \mathbf{r}(\mathbf{x}_j)$, $\mathbf{J}_{m_j} \approx \mathbf{J}(\mathbf{x}_j)$ and
$$\mathbf{s}_j \approx \operatorname{argmin}_{\mathbf{s}} m_j(\mathbf{s}) = \frac{1}{2} \|\mathbf{r}_{m_j} + \mathbf{J}_{m_j} \mathbf{s}\|^2 + \frac{\gamma_j \|\mathbf{J}_{m_j}^\top \mathbf{r}(\mathbf{x}_j)\|}{2} \|\mathbf{s}\|^2.$$
- Compute $\mathbf{r}_j^0 \approx \mathbf{r}(\mathbf{x}_j)$ and $\mathbf{r}_j^s \approx \mathbf{r}(\mathbf{x}_j + \mathbf{s}_j)$.
- If $\frac{\frac{1}{2} \|\mathbf{r}_j^0\|^2 - \frac{1}{2} \|\mathbf{r}_j^s\|^2}{m_j(0) - m_j(s)} \geq \eta$, set $\mathbf{x}_{j+1} = \mathbf{x}_j + \mathbf{s}_j$ and $\gamma_{j+1} = 0.5\gamma_j$;
- Otherwise, set $\mathbf{x}_{j+1} = \mathbf{x}_j$ and $\gamma_{j+1} = 2\gamma_j$.

- Two sources of randomness (models/estimates);
- Accounted for via martingale-type properties.

Probabilistic models

Accuracy property

For any realization, the model m_j is called $(\kappa_{ef}, \kappa_{eg})$ -accurate if

$$\|\mathbf{J}(\mathbf{x}_j)^\top \mathbf{r}(\mathbf{x}_j) - \mathbf{J}_{m_j}^\top \mathbf{r}_{m_j}\| \leq \frac{\kappa_{eg}}{\gamma_j}, \quad \left| \frac{1}{2} \|\mathbf{r}(\mathbf{x}_j)\|^2 - \frac{1}{2} \|\mathbf{r}_{m_j}\|^2 \right| \leq \frac{\kappa_{ef}}{\gamma_j^2}$$

Probabilistic accuracy property

The **random** model sequence $\{m_j\}$ is called $(p, \kappa_{ef}, \kappa_{eg})$ -accurate if

$$\forall j, \quad \mathbb{P}(m_j \text{ } (\kappa_{ef}, \kappa_{eg})\text{-accurate} \mid \mathcal{F}_{j-1}) \geq p.$$

- $\mathcal{F}_{j-1} = \sigma(m_0, \dots, m_{j-1}, \mathbf{r}_0^0, \mathbf{r}_0^s, \dots, \mathbf{r}_{j-1}^0, \mathbf{r}_{j-1}^s)$ represents the history of the algorithm up to iteration j .

Probabilistic function estimates

Accurate function estimates

Estimates $(\mathbf{r}_j^0, \mathbf{r}_j^s)$ are ε_f -accurate if

$$\left| \frac{1}{2} \|\mathbf{r}_j^0\|^2 - \frac{1}{2} \|\mathbf{r}(\mathbf{x}_j)\|^2 \right| \leq \frac{\varepsilon_f}{\gamma_j^2} \quad \text{and} \quad \left| \frac{1}{2} \|\mathbf{r}_j^s\|^2 - \frac{1}{2} \|\mathbf{r}(\mathbf{x}_j + \mathbf{s}_j)\|^2 \right| \leq \frac{\varepsilon_f}{\gamma_j^2}$$

Probabilistically accurate estimates

The random estimate sequence $\{(\mathbf{r}_j^0, \mathbf{r}_j^s)\}$ is (q, ε_f) -accurate if

$$\forall j, \quad \mathbb{P} \left((\mathbf{r}_j^0, \mathbf{r}_j^s) \text{ } \varepsilon_f\text{-accurate} \mid \mathcal{F}_{j-1/2} \right) \geq q.$$

- $\mathcal{F}_{j-1/2} = \sigma(m_0, \dots, m_{j-1}, \mathbf{m}_j, \mathbf{r}_0^0, \mathbf{r}_0^s, \dots, \mathbf{r}_{j-1}^0, \mathbf{r}_{j-1}^s)$ represents the iteration of the algorithm up to the computation of \mathbf{r}_j^0 and \mathbf{r}_j^s .

Probabilistic complexity analysis

Goal: Bound the **stopping time** $T_\epsilon = \min\{j \mid \|\mathbf{J}(\mathbf{x}_j)^\top \mathbf{r}(\mathbf{x}_j)\| \leq \epsilon_d\}$.

Theorem (BDKR '18)

If $\{m_j\}$ is $(p, \kappa_{ef}, \kappa_{eg})$ -accurate and $\{(\mathbf{r}_j^0, \mathbf{r}_j^s)\}$ is (q, ε_f) -accurate, then

$$\mathbb{E}[T_\epsilon] \leq \mathcal{O}\left(\frac{pq}{pq - 1/2} \kappa^2 \epsilon_d^{-2}\right),$$

$$\kappa = \max\{\kappa_{ef}, \kappa_{eg}, \varepsilon_f\}.$$

Probabilistic accuracy: An example

Theorem (BDKR '20)

If for every iteration j , B is approximated using an ensemble of

- $n_m = \mathcal{O} \left(\max \left\{ \frac{\|\mathbf{x}_j - \mathbf{z}_b\|_{B^{-1}}^2}{\kappa_{ef}}, \frac{\|\mathbf{B}^{-1}(\mathbf{x}_j - \mathbf{z}_b)\|}{\kappa_{eg}} \right\} n \right)$ samples to build m_j ,
- $n_f = \mathcal{O} \left(\frac{\|\mathbf{x}_j - \mathbf{z}_b\|_{B^{-1}}^2}{\varepsilon_f} n \right)$ samples to compute $(\mathbf{r}_j^0, \mathbf{r}_j^s)$,

then there exists (p, q) such that $\{m_j\}$ is $(p, \kappa_{ef}, \kappa_{eg})$ -accurate, and $\{(\mathbf{r}_j^0, \mathbf{r}_j^s)\}$ is (q, ε_f) -accurate.

- p, q depend on variance terms, likely unknown;
- $p, q \rightarrow 1$ as $n_m, n_f \rightarrow \infty$.

About the complexity results

Our problem: $\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|^2$

- Used $\|\mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x})\|$ as a complexity metric;
- Oblivious to the least-square structure;
- May want to stop when residuals are small.

About the complexity results

Our problem: $\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|^2$

- Used $\|\mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x})\|$ as a complexity metric;
- Oblivious to the least-square structure;
- May want to stop when residuals are small.

Scaled gradient (Cartis, Gould, Toint '13; Gould, Rees, Scott '19)

$$g(\mathbf{x}) := \begin{cases} \frac{\mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x})}{\|\mathbf{r}(\mathbf{x})\|} & \text{if } \|\mathbf{r}(\mathbf{x})\| > 0 \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

- Stopping criterion for complexity:

$$\|\mathbf{r}(\mathbf{x})\| \leq \epsilon_p \quad \text{or} \quad \|g(\mathbf{x})\| \leq \epsilon_d.$$

- Complexity of LM: for any $i \in \mathbb{N} \cup \{-1\}$, $\mathcal{O}(\epsilon_d^{-2} \epsilon_p^{-1/2^i})$.

New scaled gradient

Given $i \in \mathbb{N} \cup \{-1\}$,

$$\mathbf{g}^i(\mathbf{x}) := \begin{cases} \frac{\|\mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x})\|}{\|\mathbf{r}(\mathbf{x})\|^{2-2-i}} & \text{if } \|\mathbf{r}(\mathbf{x})\| \neq 0, \\ 0 & \text{otherwise.} \end{cases}$$

- Stopping criterion for complexity:

$$\|\mathbf{r}(\mathbf{x})\| \leq \epsilon_p \quad \text{or} \quad \|\mathbf{g}^i(\mathbf{x})\| \leq \epsilon_d.$$

- $i = -1$: Classical gradient;
- $i = 0$: CGT scaled gradient;
- $i \rightarrow \infty$: Resembles gradient dominance.

Complexity table

Goal: Find an (ϵ_p, ϵ_d) -point \mathbf{x}_k such that

$$\|\mathbf{r}(\mathbf{x}_k)\| \leq \epsilon_p \quad \text{or} \quad \|g^i(\mathbf{x}_k)\| \leq \epsilon_d.$$

Complexity results (BDKR '20)

i	Arbitrary	$i = -1$	$i = 0$	$i \rightarrow \infty$
$g^i(\mathbf{x})$	$\frac{\ \mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x})\ }{\ \mathbf{r}(\mathbf{x})\ ^{2-2^{-i}}}$	$\ \mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x})\ $	$\frac{\ \mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x})\ }{\ \mathbf{r}(\mathbf{x})\ }$	$\frac{\ \mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x})\ }{\ \mathbf{r}(\mathbf{x})\ ^2}$
Order	$\epsilon_d^{-2} \epsilon_p^{-(4+2^{1-i})}$	ϵ_d^{-2}	$\epsilon_d^{-2} \epsilon_p^{-2}$	$\epsilon_d^{-2} \epsilon_p^{-4}.$

- Probabilistic counterparts (expected stopping time);
- For the connoisseur: other analyzes get better bounds in terms of ϵ_p/ϵ_d but have exponential dependencies in 2^i .

Our contributions

- Redefined probabilistic property for Levenberg-Marquardt schemes (analogy with trust region);
- Complexity analysis for **stochastic function estimates**;
- Application in a data assimilation setting.
- A family of **complexity metrics** and results.

Derivative-free Levenberg-Marquardt

A stochastic Levenberg-Marquardt method using random models with complexity results and application to data assimilation.

E. Bergou, Y. Diouane, V. Kungurstev and C. W. Royer.

<https://arXiv.org/abs/1807.02176v2>

A nonmonotone matrix-free algorithm for nonlinear equality-constrained least-squares problems.

E. Bergou, Y. Diouane, V. Kungurstev and C. W. Royer.

<https://arXiv.org/abs/arXiv:2006.16340v2>.

- Addressed (among others) the original 4Dvar formulation:

$$\min_{\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_T \in \mathbb{R}^n} \quad \frac{1}{2} \|\mathbf{z}_0 - \mathbf{z}_b\|_{\mathbf{B}^{-1}}^2 + \frac{1}{2} \sum_{i=0}^T \|\mathbf{y}_i - \mathcal{H}(\mathbf{z}_i)\|_{\mathbf{R}_i^{-1}}^2$$

s.t. $\mathbf{z}_i = \mathcal{M}_i(\mathbf{z}_{i-1}), \quad i = 1, \dots, T.$

- Derivative-based approach but **inexact**.

Going forward (CSE 23?)

- Combine random models with constraints;
- Practical probabilistic approaches.