Algorithms and application for special classes of nonlinear least squares problems 2023

Clément W. Royer

US-Mexico workshop on optimization and its applications

## January 9, 2023



Nonlinear Least Squares Problems 2023

## Not so long ago, not so far away...

BULL. AUSTRAL. MATH. SOC.	90C30
VOL. 31 (1985), 309-311.	(65K05)

ALGORITHMS AND APPLICATION FOR SPECIAL CLASSES OF NONLINEAR LEAST SQUARES PROBLEMS

STEPHEN J. WRIGHT

A Nonlinear Least Squares problem is an optimization problem for which the objective function to be minimized has the form

$$F(x) = \sum_{i=1}^m f_i^2(x)$$
 ,  $x \in \operatorname{R}^n$  ,  $m \ge n$  .

- Nonlinear least squares + constraints/nonsmoothness;
- Efficient solve for sparse Jacobian;
- Cool application (geophysics).

What more can I do?

# Rebooting the (least-squares) franchise



A DIONAL LINEAR THE ADVISOR DEVICE TO THE ADVISOR OF THE ADVISOR O

### Nonlinear least squares

- A basic problem...
- ...with modern instances.

## Revisit algorithms

- Complexity bounds for Gauss-Newton methods.
- Complexity metrics.

## Go a bit further

- Inexactness and stochasticity.
- Probabilistic results.

- Problem and first results
- 2 More complexity results
- Beyond the deterministic setting
- Application: Learning dynamics

## Problem and first results

- 2 More complexity results
- 3 Beyond the deterministic setting
- 4 Application: Learning dynamics

## General setup

## Nonlinear least-squares

$$\min_{\boldsymbol{x}\in\mathbb{R}^n} f(\boldsymbol{x}) := \frac{1}{2} \|\boldsymbol{r}(\boldsymbol{x})\|^2, \qquad \boldsymbol{r}:\mathbb{R}^n \mapsto \mathbb{R}^m, \boldsymbol{r}\in\mathcal{C}^2.$$
Jacobian:  $\boldsymbol{J}(\boldsymbol{x}) := [\nabla r_i(\boldsymbol{x})^\top] \in \mathbb{R}^{m \times n}.$ 

### Nonlinear least-squares

$$\min_{oldsymbol{x}\in\mathbb{R}^n}f(oldsymbol{x}):=rac{1}{2}\|oldsymbol{r}(oldsymbol{x})\|^2,\qquadoldsymbol{r}:\mathbb{R}^n\mapsto\mathbb{R}^m,oldsymbol{r}\in\mathcal{C}^2.$$

**Jacobian:**  $J(\mathbf{x}) := [\nabla r_i(\mathbf{x})^\top] \in \mathbb{R}^{m \times n}$ .

## Gauss-Newton techniques

- Gauss-Newton model  $f(\mathbf{x} + \mathbf{s}) \approx \frac{1}{2} \|\mathbf{r}(\mathbf{x}) + \mathbf{J}(\mathbf{x})\mathbf{s}\|^2$ ;
- Steps computed via line search/trust region;
- Hessian approximated by  $J(x)^{T}J(x)$ .

### Nonlinear least-squares

$$\min_{oldsymbol{x}\in\mathbb{R}^n}f(oldsymbol{x}):=rac{1}{2}\|oldsymbol{r}(oldsymbol{x})\|^2,\qquadoldsymbol{r}:\mathbb{R}^n\mapsto\mathbb{R}^m,oldsymbol{r}\in\mathcal{C}^2.$$

**Jacobian:**  $J(\mathbf{x}) := [\nabla r_i(\mathbf{x})^\top] \in \mathbb{R}^{m \times n}$ .

### Gauss-Newton techniques

- Gauss-Newton model  $f(\mathbf{x} + \mathbf{s}) \approx \frac{1}{2} \|\mathbf{r}(\mathbf{x}) + \mathbf{J}(\mathbf{x})\mathbf{s}\|^2$ ;
- Steps computed via line search/trust region;
- Hessian approximated by  $J(x)^{T}J(x)$ .

## Levenberg(-Morrison)-Marquardt

- Regularized 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  $\gamma$  set adaptively.

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

Inputs:  $\mathbf{x}_0 \in \mathbb{R}^n, \gamma_0 \ge \gamma_{\min} > 0, \eta > 0.$ Iteration k: Given  $(\mathbf{x}_k, \gamma_k)$ ,

Compute

$$oldsymbol{s}_k pprox rgmin_{oldsymbol{s}} m_k(oldsymbol{s}) := rac{1}{2} \|oldsymbol{r}(oldsymbol{x}_k) + oldsymbol{J}(oldsymbol{x}_k) oldsymbol{s} \|^2 + rac{\gamma_k}{2} \|oldsymbol{s}\|^2.$$

$$\frac{\frac{1}{2}\|r(\boldsymbol{x}_k)\|^2 - \frac{1}{2}\|\boldsymbol{r}(\boldsymbol{x}_k + \boldsymbol{s}_k)\|^2}{m_k(0) - m_k(\boldsymbol{s}_k)} \geq \eta,$$

set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{s}_k$  and  $\gamma_{k+1} = \max\{0.5\gamma_k, \gamma_{\min}\};$ 

• Otherwise, set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k$  and  $\gamma_{k+1} = 2\gamma_k$ .

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

Inputs:  $\mathbf{x}_0 \in \mathbb{R}^n, \gamma_0 \ge \gamma_{\min} > 0, \eta > 0.$ Iteration k: Given  $(\mathbf{x}_k, \gamma_k)$ ,

Compute

$$oldsymbol{s}_k pprox rgmin_{oldsymbol{s}} m_k(oldsymbol{s}) := rac{1}{2} \|oldsymbol{r}(oldsymbol{x}_k) + oldsymbol{J}(oldsymbol{x}_k) oldsymbol{s} \|^2 + rac{\gamma_k}{2} \|oldsymbol{s}\|^2.$$

$$\frac{\frac{1}{2}\|r(\boldsymbol{x}_k)\|^2 - \frac{1}{2}\|r(\boldsymbol{x}_k + \boldsymbol{s}_k)\|^2}{m_k(0) - m_k(\boldsymbol{s}_k)} \geq \eta_{2}$$

set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{s}_k$  and  $\gamma_{k+1} = \max\{0.5\gamma_k, \gamma_{\min}\};$ 

• Otherwise, set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k$  and  $\gamma_{k+1} = 2\gamma_k$ .

Goal: Prove a complexity result for the method.

## Complexity in nonconvex optimization

**Setup:** Sequence of points  $\{x_k\}$  generated by an algorithm applied to  $\min_{x \in \mathbb{R}^n} f(x)$ .

## Complexity in nonconvex optimization

**Setup:** Sequence of points  $\{x_k\}$  generated by an algorithm applied to  $\min_{x \in \mathbb{R}^n} f(x)$ .

#### Standard complexity result

Given  $\epsilon_g \in (0, 1)$ :

- Worst-case number of iterations to obtain  $\boldsymbol{x}_k$  such that  $\|\nabla f(\boldsymbol{x}_k)\| \leq \epsilon_g$ .
- Focus: Dependency on  $\epsilon_g$ .

# Complexity in nonconvex optimization

**Setup:** Sequence of points  $\{x_k\}$  generated by an algorithm applied to  $\min_{x \in \mathbb{R}^n} f(x)$ .

### Standard complexity result

Given  $\epsilon_g \in (0, 1)$ :

- Worst-case number of iterations to obtain  $\boldsymbol{x}_k$  such that  $\|\nabla f(\boldsymbol{x}_k)\| \leq \epsilon_g$ .
- Focus: Dependency on  $\epsilon_g$ .

#### Some examples

- Gradient descent:  $\mathcal{O}(\epsilon_g^{-2})$  iterations.
- Newton:  $\mathcal{O}(\epsilon_g^{-2})$  iterations.
- Cubic regularization/Modified Newton:  $\mathcal{O}(\epsilon_g^{-3/2})$  iterations.

- Problem:  $\min_{\boldsymbol{x} \in \mathbb{R}^n} \frac{1}{2} \|\boldsymbol{r}(\boldsymbol{x})\|^2$ :
- Goal: Find  $\boldsymbol{x}_k$  such that  $\|\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x}_k)\| \leq \epsilon_g$ .

- Problem:  $\min_{\boldsymbol{x} \in \mathbb{R}^n} \frac{1}{2} \|\boldsymbol{r}(\boldsymbol{x})\|^2$ :
- Goal: Find  $\boldsymbol{x}_k$  such that  $\|\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x}_k)\| \leq \epsilon_g$ .

#### Some results

- Gradient descent:  $\mathcal{O}(\epsilon_g^{-2})$  iterations.
- Gauss-Newton + line search/trust region:  $\mathcal{O}(\epsilon_g^{-2})$  iterations.
- Levenberg-Marquardt:  $\mathcal{O}(\epsilon_g^{-2})$  iterations.

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

Inputs:  $\mathbf{x}_0 \in \mathbb{R}^n, \gamma_0 \ge \gamma_{\min} > 0, \eta > 0.$ Iteration k: Given  $(\mathbf{x}_k, \gamma_k)$ ,

Compute

$$oldsymbol{s}_k pprox rgmin_{oldsymbol{s}} m_k(oldsymbol{s}) := rac{1}{2} \|oldsymbol{r}(oldsymbol{x}_k) + oldsymbol{J}(oldsymbol{x}_k) oldsymbol{s} \|^2 + rac{\gamma_k}{2} \|oldsymbol{s}\|^2.$$

$$\frac{\frac{1}{2}\|r(\boldsymbol{x}_k)\|^2 - \frac{1}{2}\|\boldsymbol{r}(\boldsymbol{x}_k + \boldsymbol{s}_k)\|^2}{m_k(0) - m_k(\boldsymbol{s}_k)} \geq \eta,$$

set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{s}_k$  and  $\gamma_{k+1} = \max\{0.5\gamma_k, \gamma_{\min}\};$ 

• Otherwise, set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k$  and  $\gamma_{k+1} = 2\gamma_k$ .

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

Inputs:  $\mathbf{x}_0 \in \mathbb{R}^n, \gamma_0 \ge \gamma_{\min} > 0, \eta > 0.$ Iteration k: Given  $(\mathbf{x}_k, \gamma_k)$ ,

Compute

$$oldsymbol{s}_k pprox rgmin_{oldsymbol{s}} m_k(oldsymbol{s}) := rac{1}{2} \|oldsymbol{r}(oldsymbol{x}_k) + oldsymbol{J}(oldsymbol{x}_k) oldsymbol{s} \|^2 + rac{\gamma_k}{2} \|oldsymbol{s}\|^2.$$

$$\frac{\frac{1}{2}\|r(\boldsymbol{x}_k)\|^2 - \frac{1}{2}\|\boldsymbol{r}(\boldsymbol{x}_k + \boldsymbol{s}_k)\|^2}{m_k(0) - m_k(\boldsymbol{s}_k)} \geq \eta,$$

set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{s}_k$  and  $\gamma_{k+1} = \max\{0.5\gamma_k, \gamma_{\min}\};$ 

• Otherwise, set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k$  and  $\gamma_{k+1} = 2\gamma_k$ .

Goal: Prove a complexity result for the method.

## Decrease guarantee

For any successful iteration  $(\mathbf{x}_{k+1} \neq \mathbf{x}_k)$ ,

$$\|r(\boldsymbol{x}_k)\|^2 - \|r(\boldsymbol{x}_{k+1})\|^2 \ge \mathcal{O}\left(\frac{\|\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x}_k)\|^2}{\gamma_k}\right).$$

#### Decrease guarantee

For any successful iteration  $(\mathbf{x}_{k+1} \neq \mathbf{x}_k)$ ,

$$\|r(\boldsymbol{x}_k)\|^2 - \|r(\boldsymbol{x}_{k+1})\|^2 \ge \mathcal{O}\left(\frac{\|\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x}_k)\|^2}{\gamma_k}\right)$$

#### Regularization parameter

• If  $\gamma_k$  large enough, the iteration is successful.

• 
$$\gamma_k \leq \gamma_{\max}$$
 for all  $k$ .

#### Decrease guarantee

For any successful iteration  $(\mathbf{x}_{k+1} \neq \mathbf{x}_k)$ ,

$$\|r(\boldsymbol{x}_k)\|^2 - \|r(\boldsymbol{x}_{k+1})\|^2 \ge \mathcal{O}\left(\frac{\|\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x}_k)\|^2}{\gamma_k}\right)$$

### Regularization parameter

• If  $\gamma_k$  large enough, the iteration is successful.

• 
$$\gamma_k \leq \gamma_{\max}$$
 for all  $k$ .

## Complexity of standard Levenberg-Marquardt

The method reaches  $\boldsymbol{x}_k$  such that  $\|\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x}_k)\| \leq \epsilon_g$  in at most  $\mathcal{O}(\epsilon_g^{-2})$  iterations.

## Problem and first results

- 2 More complexity results
  - 3 Beyond the deterministic setting
  - 4 Application: Learning dynamics

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

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

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

- Used  $\| \boldsymbol{J}(\boldsymbol{x})^{\mathrm{T}} \boldsymbol{r}(\boldsymbol{x}) \|$  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)

$$oldsymbol{g}(oldsymbol{x}) := \left\{egin{array}{cc} rac{oldsymbol{J}(oldsymbol{x})^{ ext{T}}oldsymbol{r}(oldsymbol{x})}{\|oldsymbol{r}(oldsymbol{x})\|} & ext{if } \|oldsymbol{r}(oldsymbol{x})\| > 0 \\ 0 & ext{otherwise.} \end{array}
ight.$$

• Stopping criterion for complexity:

$$\|\boldsymbol{r}(\boldsymbol{x})\| \leq \epsilon_r \quad \text{or} \quad \|\boldsymbol{g}(\boldsymbol{x})\| \leq \epsilon_g.$$

# An alternate complexity measure ('ed)

**Goal:** Find  $x_k$  such that

$$\|\boldsymbol{r}(\boldsymbol{x}_k)\| \leq \epsilon_r \quad \text{or} \quad \|\boldsymbol{g}(\boldsymbol{x}_k)\| \leq \epsilon_g, \quad \boldsymbol{g}(\boldsymbol{x}_k) := \begin{cases} rac{\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}} \boldsymbol{r}(\boldsymbol{x}_k)}{\|\boldsymbol{r}(\boldsymbol{x}_k)\|} & ext{if } \|\boldsymbol{r}(\boldsymbol{x}_k)\| > 0\\ 0 & ext{otherwise.} \end{cases}$$

**Goal:** Find  $\boldsymbol{x}_k$  such that

$$\|m{r}(m{x}_k)\| \leq \epsilon_r \quad ext{or} \quad \|m{g}(m{x}_k)\| \leq \epsilon_g, \quad m{g}(m{x}_k) := \left\{egin{array}{c} rac{m{J}(m{x}_k)^{ op}m{r}(m{x}_k)}{\|m{r}(m{x}_k)\|} & ext{if} \ \|m{r}(m{x}_k)\| > 0 \ 0 & ext{otherwise.} \end{array}
ight.$$

Complexity of Levenberg-Marquardt (Gould, Rees, Scott '19)

For any  $i \in \mathbb{N} \cup \{-1\}$ , the method finds a suitable  $\boldsymbol{x}_k$  in at most

 $\mathcal{O}(2^i \epsilon_g^{-2} \epsilon_r^{-1/2^i})$  iterations.

**Goal:** Find  $\boldsymbol{x}_k$  such that

$$\|\boldsymbol{r}(\boldsymbol{x}_k)\| \leq \epsilon_r \quad ext{or} \quad \|\boldsymbol{g}(\boldsymbol{x}_k)\| \leq \epsilon_g, \quad \boldsymbol{g}(\boldsymbol{x}_k) := \left\{ egin{array}{c} rac{m{J}(m{x}_k)^{ ext{T}}m{r}(m{x}_k)}{\|m{r}(m{x}_k)\|} & ext{if } \|m{r}(m{x}_k)\| > 0 \\ 0 & ext{otherwise.} \end{array} 
ight.$$

Complexity of Levenberg-Marquardt (Gould, Rees, Scott '19)

For any  $i \in \mathbb{N} \cup \{-1\}$ , the method finds a suitable  $\boldsymbol{x}_k$  in at most

 $\mathcal{O}(2^i \epsilon_g^{-2} \epsilon_r^{-1/2^i})$  iterations.

• Part of more results on high-order regularization methods.

• Asymptotically: 
$$\epsilon_r^{-1/2'} \to 1$$
 but  $2^i \to \infty$ .

# Alternate metric (Bergou, Diouane, Kungurstev, R '22)

#### New scaled gradient

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

$$oldsymbol{g}^{\mathrm{i}}(oldsymbol{x}) \ := \ \left\{ egin{array}{cc} rac{\|oldsymbol{J}(oldsymbol{x})^{ op}oldsymbol{r}(oldsymbol{x})\|}{\|oldsymbol{r}(oldsymbol{x})\|^{2-2^{-\mathrm{i}}}} & \mathrm{if} \ \|oldsymbol{r}(oldsymbol{x})\| 
eq 0, \ 0 & \mathrm{otherwise.} \end{array} 
ight.$$

• Stopping criterion for complexity:

$$\|\boldsymbol{r}(\boldsymbol{x})\| \leq \epsilon_r \quad \text{or} \quad \|\boldsymbol{g}^{\mathrm{i}}(\boldsymbol{x})\| \leq \epsilon_g.$$

# Alternate metric (Bergou, Diouane, Kungurstev, R '22)

### New scaled gradient

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

$$oldsymbol{g}^{\mathrm{i}}(oldsymbol{x}) \ := \ \left\{ egin{array}{cc} rac{\|oldsymbol{J}(oldsymbol{x})^{ op}oldsymbol{r}(oldsymbol{x})\|}{\|oldsymbol{r}(oldsymbol{x})\|^{2-2^{-\mathrm{i}}}} & \mathrm{if} \ \|oldsymbol{r}(oldsymbol{x})\| 
eq 0, \ 0 & \mathrm{otherwise.} \end{array} 
ight.$$

• Stopping criterion for complexity:

$$\|\boldsymbol{r}(\boldsymbol{x})\| \leq \epsilon_r \quad \text{or} \quad \|\boldsymbol{g}^{\mathrm{i}}(\boldsymbol{x})\| \leq \epsilon_g.$$

- i = -1: Classical gradient;
- i = 0: CGT scaled gradient;
- $\mathfrak{i} \to \infty$ :  $\|\boldsymbol{g}^{\mathfrak{i}}(\boldsymbol{x})\| > \epsilon_{g}$  akin to gradient dominance.

#### Decrease guarantees

For any successful iteration  $(\mathbf{x}_{k+1} \neq \mathbf{x}_k)$ ,

$$\|\boldsymbol{r}(\boldsymbol{x}_k)\|^2 - \|\boldsymbol{r}(\boldsymbol{x}_{k+1})\|^2 \ge \mathcal{O}\left(\frac{\|\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x}_k)\|^2}{\gamma_k}\right)$$

and (if  $\|\boldsymbol{r}(\boldsymbol{x}_k)\| \neq 0$ )

$$\|m{r}(m{x}_k)\|^{rac{1}{2^{\mathfrak{i}}}} - \|m{r}(m{x}_{k+1})\|^{rac{1}{2^{\mathfrak{i}}}} \geq \mathcal{O}\left(rac{\|m{g}^{\mathfrak{i}}(m{x}_k)\|^2\|m{r}(m{x}_k)\|^{(4-2^{1-\mathfrak{i}})}}{\gamma_k}
ight).$$

#### Decrease guarantees

For any successful iteration  $(\mathbf{x}_{k+1} \neq \mathbf{x}_k)$ ,

$$\|\boldsymbol{r}(\boldsymbol{x}_k)\|^2 - \|\boldsymbol{r}(\boldsymbol{x}_{k+1})\|^2 \ge \mathcal{O}\left(\frac{\|\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x}_k)\|^2}{\gamma_k}\right)$$

and (if  $\|\boldsymbol{r}(\boldsymbol{x}_k)\| \neq 0$ )

$$\|\boldsymbol{r}(\boldsymbol{x}_{k})\|^{\frac{1}{2^{i}}} - \|\boldsymbol{r}(\boldsymbol{x}_{k+1})\|^{\frac{1}{2^{i}}} \ge \mathcal{O}\left(\frac{\|\boldsymbol{g}^{i}(\boldsymbol{x}_{k})\|^{2}\|\boldsymbol{r}(\boldsymbol{x}_{k})\|^{(4-2^{1-i})}}{\gamma_{k}}\right)$$

### Regularization parameter

- If  $\gamma_k$  large enough, the iteration is successful.
- $\gamma_k \leq \gamma_{\max}$  for all k.

# Complexity table

## **Goal:** Find $\boldsymbol{x}_k$ such that

$$\|\boldsymbol{r}(\boldsymbol{x}_k)\| \leq \epsilon_r \quad \text{or} \quad \|\boldsymbol{g}^{\mathrm{i}}(\boldsymbol{x}_k)\| \leq \epsilon_g.$$

## Complexity results (BDKR '22)

i	Arbitrary	$\mathfrak{i} = -1$	i = 0
$g^{i}(\mathbf{x})$	$\frac{\ \boldsymbol{J}(\boldsymbol{x})^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x})\ }{\ \boldsymbol{r}(\boldsymbol{x})\ ^{2-2^{-1}}}$	$\ J(x)^{\mathrm{T}}r(x)\ $	$\frac{\ J(x)^{\mathrm{T}}r(x)\ }{\ r(x)\ }$
Order	$\left  \epsilon_g^{-2} \epsilon_r^{-(4-2^{1-i})} \right $	$\epsilon_g^{-2}$	$\epsilon_g^{-2}\epsilon_r^{-2}$ .

- Matches existing results for i = -1.
- For i = 0: previous results get better bounds in terms of ε<sub>r</sub> but with very large constants (2<sup>i</sup>).

- Problem and first results
- 2 More complexity results
- Beyond the deterministic setting
  - 4 Application: Learning dynamics

Stochastic nonlinear least-squares

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

• Values of *r* and Jacobian *J* only accessed through stochastic estimates.

## Challenges

- Every evaluation is replaced by a random estimate;
- Decrease no longer guaranteed;
- Accuracy of evaluation matters.

## Using inexact values

Inputs:  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\gamma_0 \ge \gamma_{\min} > 0$ . Iteration k: Given  $(\mathbf{x}_k, \gamma_k)$ ,

- Compute  $\mathbf{r}_{m_k} \approx \mathbf{r}(\mathbf{x}_k)$ ,  $\mathbf{J}_{m_k} \approx \mathbf{J}(\mathbf{x}_k)$  and  $\mathbf{s}_k \approx \operatorname{argmin}_{\mathbf{s}} m_k(\mathbf{s}) = \frac{1}{2} \|\mathbf{r}_{m_k} + \mathbf{J}_{m_k} \mathbf{s}\|^2 + \frac{\gamma_k}{2} \|\mathbf{s}\|^2$ .
- Compute  $\mathbf{r}_k^0 \approx \mathbf{r}(\mathbf{x}_k)$  and  $\mathbf{r}_k^s \approx \mathbf{r}(\mathbf{x}_k + \mathbf{s}_k)$ .
- If  $\frac{\frac{1}{2} \|\boldsymbol{r}_{k}^{0}\|^{2} \frac{1}{2} \|\boldsymbol{r}_{k}^{s}\|^{2}}{m_{k}(0) m_{k}(s)} \geq \eta$ , set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_{k} + \boldsymbol{s}_{k}$  and  $\gamma_{k+1} = \max\{0.5\gamma_{k}, \gamma_{\min}\};$
- Otherwise, set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k$  and  $\gamma_{k+1} = 2\gamma_k$ .

## Using inexact values

Inputs:  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\gamma_0 \ge \gamma_{\min} > 0$ . Iteration k: Given  $(\mathbf{x}_k, \gamma_k)$ ,

- Compute  $\mathbf{r}_{m_k} \approx \mathbf{r}(\mathbf{x}_k)$ ,  $\mathbf{J}_{m_k} \approx \mathbf{J}(\mathbf{x}_k)$  and  $\mathbf{s}_k \approx \operatorname{argmin}_{\mathbf{s}} m_k(\mathbf{s}) = \frac{1}{2} \|\mathbf{r}_{m_k} + \mathbf{J}_{m_k} \mathbf{s}\|^2 + \frac{\gamma_k}{2} \|\mathbf{s}\|^2$ .
- Compute  $\mathbf{r}_k^0 \approx \mathbf{r}(\mathbf{x}_k)$  and  $\mathbf{r}_k^s \approx \mathbf{r}(\mathbf{x}_k + \mathbf{s}_k)$ .
- If  $\frac{\frac{1}{2} \|r_k^0\|^2 \frac{1}{2} \|r_k^s\|^2}{m_k(0) m_k(s)} \ge \eta$ , set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k + \boldsymbol{s}_k$  and  $\gamma_{k+1} = \max\{0.5\gamma_k, \gamma_{\min}\};$
- Otherwise, set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k$  and  $\gamma_{k+1} = 2\gamma_k$ .

**Goal:** Prove a complexity result for this inexact method. **Key:** Use  $\gamma_k$  to monitor the inexactness and the convergence.

# Complexity analysis in an inexact setting

## Accuracy requirements (model)

For every k,

$$\| oldsymbol{J}(oldsymbol{x}_k)^{ op}oldsymbol{r}(oldsymbol{x}_k) - oldsymbol{J}_{m_k}^{ op}oldsymbol{r}_{m_k} \| \leq \mathcal{O}\left(rac{1}{\gamma_k}
ight)$$

and

$$\left|rac{1}{2}\|oldsymbol{r}(x_k)\|^2 - rac{1}{2}\|oldsymbol{r}_{m_k}\|^2
ight| \leq \mathcal{O}\left(rac{1}{\gamma_k}
ight).$$

## Accuracy requirements (evaluations)

For every k,

$$\left\| \frac{1}{2} \| \boldsymbol{r}_k^0 \|^2 - \frac{1}{2} \| \boldsymbol{r}(\boldsymbol{x}_k) \|^2 \right\| \leq \mathcal{O}\left( \frac{1}{\gamma_k^2} \right)$$

and

$$\left| rac{1}{2} \| oldsymbol{r}_k^{oldsymbol{s}} \|^2 - rac{1}{2} \| oldsymbol{r}(oldsymbol{x}_k + oldsymbol{s}_k) \|^2 
ight| \leq \mathcal{O}\left(rac{1}{\gamma_k^2}
ight)$$

# Complexity analysis

#### Using inexactness

• With the same theory as in the exact case, get  $\mathcal{O}(\epsilon^{-3})$  instead of  $\mathcal{O}(\epsilon^{-2})$  to obtain  $\|J(\mathbf{x}_k)^{\mathrm{T}} \mathbf{r}(\mathbf{x}_k)\| \leq \epsilon!$ 

# Complexity analysis

#### Using inexactness

- With the same theory as in the exact case, get  $\mathcal{O}(\epsilon^{-3})$  instead of  $\mathcal{O}(\epsilon^{-2})$  to obtain  $\|\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x}_k)\| \leq \epsilon!$
- Arguments:
  - Still decrease in  $\mathcal{O}\left(\frac{\|\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x}_k)\|^2}{\gamma_k}\right);$
  - But now  $\gamma_k$  grows as  $\mathcal{O}(\epsilon^{-1})!$

# Complexity analysis

#### Using inexactness

- With the same theory as in the exact case, get  $\mathcal{O}(\epsilon^{-3})$  instead of  $\mathcal{O}(\epsilon^{-2})$  to obtain  $\|\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x}_k)\| \leq \epsilon!$
- Arguments:
  - Still decrease in  $\mathcal{O}\left(\frac{\|\boldsymbol{J}(\boldsymbol{x}_k)^{\mathrm{T}}\boldsymbol{r}(\boldsymbol{x}_k)\|^2}{\gamma_k}\right);$
  - But now  $\gamma_k$  grows as  $\mathcal{O}(\epsilon^{-1})!$

## A fix (BDKR '22)

- The analysis reveals  $\gamma = \mathcal{O}(\|J(\mathbf{x})^{\top}r(\mathbf{x})\|/\|\mathbf{s}\|);$
- By analogy with trust-region, we want  $\gamma = 1/\| \pmb{s} \|;$
- A scaling will help us achieve that.

## A corrected Levenberg-Marquardt method

Inputs:  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\gamma_0 \ge \gamma_{\min} > 0$ . Iteration k: Given  $(\mathbf{x}_k, \gamma_k)$ ,

• Compute  $\mathbf{r}_{m_k} \approx \mathbf{r}(\mathbf{x}_k), \ \mathbf{J}_{m_k} \approx \mathbf{J}(\mathbf{x}_k)$  and

$$oldsymbol{s}_k pprox ext{argmin}_{oldsymbol{s}} m_k(oldsymbol{s}) = rac{1}{2} \|oldsymbol{r}_{m_k} + oldsymbol{J}_{m_k} oldsymbol{s}\|^2 + rac{\gamma_k \|oldsymbol{J}_{m_k} oldsymbol{r}_{m_k}\|}{2} \|oldsymbol{s}\|^2.$$

- Compute  $\mathbf{r}_k^0 \approx \mathbf{r}(\mathbf{x}_k)$  and  $\mathbf{r}_k^s \approx \mathbf{r}(\mathbf{x}_k + \mathbf{s}_k)$ .
- If  $\frac{\frac{1}{2} \|\boldsymbol{r}_{k}^{0}\|^{2} \frac{1}{2} \|\boldsymbol{r}_{k}^{s}\|^{2}}{m_{k}(0) m_{k}(s)} \geq \eta$ , set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_{k} + \boldsymbol{s}_{k}$  and  $\gamma_{k+1} = \max\{0.5\gamma_{k}, \gamma_{\min}\};$
- Otherwise, set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k$  and  $\gamma_{k+1} = 2\gamma_k$ .

## A corrected Levenberg-Marquardt method

Inputs:  $\mathbf{x}_0 \in \mathbb{R}^n$ ,  $\gamma_0 \ge \gamma_{\min} > 0$ . Iteration k: Given  $(\mathbf{x}_k, \gamma_k)$ ,

• Compute  $\mathbf{r}_{m_k} \approx \mathbf{r}(\mathbf{x}_k), \ \mathbf{J}_{m_k} \approx \mathbf{J}(\mathbf{x}_k)$  and

 $\boldsymbol{s}_k \approx \operatorname{argmin}_{\boldsymbol{s}} m_k(\boldsymbol{s}) = \frac{1}{2} \|\boldsymbol{r}_{m_k} + \boldsymbol{J}_{m_k} \boldsymbol{s}\|^2 + \frac{\gamma_k \|\boldsymbol{J}_{m_k}^\top \boldsymbol{r}_{m_k}\|}{2} \|\boldsymbol{s}\|^2.$ 

- Compute  $\mathbf{r}_k^0 \approx \mathbf{r}(\mathbf{x}_k)$  and  $\mathbf{r}_k^s \approx \mathbf{r}(\mathbf{x}_k + \mathbf{s}_k)$ .
- If  $\frac{\frac{1}{2} \|\boldsymbol{r}_{k}^{0}\|^{2} \frac{1}{2} \|\boldsymbol{r}_{k}^{s}\|^{2}}{m_{k}(0) m_{k}(s)} \geq \eta$ , set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_{k} + \boldsymbol{s}_{k}$  and  $\gamma_{k+1} = \max\{0.5\gamma_{k}, \gamma_{\min}\};$
- Otherwise, set  $\boldsymbol{x}_{k+1} = \boldsymbol{x}_k$  and  $\gamma_{k+1} = 2\gamma_k$ .
- Two sources of inexactness (models/estimates);
- Analysis can be deterministic or probabilistic.

## Probabilistic models

#### Accuracy property

For any realization,  $(\boldsymbol{J}_{m_k}, \boldsymbol{r}_{m_k})$  is called accurate if

$$\|oldsymbol{J}(oldsymbol{x}_k)^{ op}oldsymbol{r}(oldsymbol{x}_k) - oldsymbol{J}_{oldsymbol{m}_k}^{ op}oldsymbol{r}_{m_k}\| \leq \mathcal{O}\left(rac{1}{\gamma_k}
ight) \quad \left|rac{1}{2}\|oldsymbol{r}(oldsymbol{x}_k)\|^2 - rac{1}{2}\|oldsymbol{r}_{m_k}\|^2
ight| \leq \mathcal{O}\left(rac{1}{\gamma_k^2}
ight)$$

#### Probabilistic accuracy property

The random model sequence  $\{(J_{m_k}, r_{m_k})\}$  is called *p*-accurate if

$$\forall k, \quad \mathbb{P}\left(\left(oldsymbol{J}_{m_k}, oldsymbol{r}_{m_k}
ight) ext{ accurate } |\mathcal{F}_{k-1}
ight) \geq p.$$

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

## Accurate function estimates

$$\begin{aligned} \left| \frac{1}{2} \| \boldsymbol{r}_k^0 \|^2 - \frac{1}{2} \| \boldsymbol{r}(\boldsymbol{x}_k) \|^2 \right| &\leq \mathcal{O}\left( \frac{1}{\gamma_k^2} \right) \\ \left| \frac{1}{2} \| \boldsymbol{r}_k^s \|^2 - \frac{1}{2} \| \boldsymbol{r}(\boldsymbol{x}_k + \boldsymbol{s}_k) \|^2 \right| &\leq \mathcal{O}\left( \frac{1}{\gamma_k^2} \right). \end{aligned}$$

### Probabilistically accurate estimates

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

$$orall k, \quad \mathbb{P}\left((\pmb{r}_k^0, \pmb{r}_k^1) ext{ accurate } ig|\mathcal{F}_{k-1/2}
ight) \geq q.$$

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

## Goal: Bound the stopping time

$$\mathcal{K}_{\epsilon} = \min\{k \; ||| oldsymbol{r}(oldsymbol{x})|| \leq \epsilon_r \quad ext{or} \quad \|oldsymbol{g}^{ ext{i}}(oldsymbol{x})\| \leq \epsilon_g\}.$$

## Goal: Bound the stopping time

$$\mathcal{K}_{\epsilon} = \min\{k \mid \| \boldsymbol{r}(\boldsymbol{x}) \| \leq \epsilon_r \quad ext{or} \quad \| \boldsymbol{g}^{ ext{i}}(\boldsymbol{x}) \| \leq \epsilon_g \}.$$

## Theorem (BDKR '22)

If  $\{(\pmb{J}_{m_k},\pmb{r}_{m_k})\}$  are *p*-accurate and  $\{(\pmb{r}_k^0,\pmb{r}_k^s)\}$  are *q*-accurate, then

$$\mathbb{E}[K_{\epsilon}] \leq \mathcal{O}\left(\frac{pq}{pq-1/2} \epsilon_g^{-2} \epsilon_r^{-(4-2^{1-i})}\right).$$

- Problem and first results
- 2 More complexity results
- 3 Beyond the deterministic setting
- Application: Learning dynamics

# Motivation: Learning ODE parameters

## Problem

• Data:  $\{z(t_i)\}_{i=0}^m$  obtained from the solution z(t) of an ODE

$$\frac{d\boldsymbol{z}}{dt}(t) = \phi_{\boldsymbol{A}}(\boldsymbol{z}(t)).$$

with  $\boldsymbol{z}(0) = \boldsymbol{z}_0$ .

• Goal: Learn the parameters  $\boldsymbol{A}$  of the dynamics  $\phi$ .

# Motivation: Learning ODE parameters

## Problem

• Data:  $\{z(t_i)\}_{i=0}^m$  obtained from the solution z(t) of an ODE

$$\frac{d\boldsymbol{z}}{dt}(t) = \phi_{\boldsymbol{A}}(\boldsymbol{z}(t)).$$

with  $z(0) = z_0$ .

• Goal: Learn the parameters  $\boldsymbol{A}$  of the dynamics  $\phi$ .

## Model: NeuralODE

• A neural network defined as the solution of an ODE:  $\pmb{z}\mapsto \pmb{y}(1)$ , where  $\pmb{y}$  solution of

$$\frac{d\boldsymbol{y}}{dt}(t) = \phi_{\boldsymbol{X}}(\boldsymbol{y}(t))$$

with  $\boldsymbol{y}(0) = \boldsymbol{y}_0$ .

• Goal: Learn X close to A.

# Illustration: Linear ODE

## Problem

- Noisy data  $\{z_i\}_{i=0}^m$  generated by a linear ODE  $\frac{dz}{dt}(t) = Az(t)$ ;
- Closed-form expression:  $z(t) = e^{At}z(0)$ .

## Training problem

$$\underset{\boldsymbol{X} \in \mathbb{R}^{n \times n}}{\text{minimize}} \frac{1}{m} \sum_{i=1}^{m} \left\| \left( \boldsymbol{I} + \frac{\boldsymbol{X}}{\ell} \right)^{\ell} \boldsymbol{z}_{i} - \boldsymbol{z}_{i+1} \right\|^{2}$$

Euler's formula:

$$e^{oldsymbol{X}} pprox \left(oldsymbol{I} + rac{oldsymbol{X}}{\ell}
ight)^\ell, \ell \geq 1.$$

 Nonconvex nonlinear least squares for ℓ ≥ 2 (strict, even high-order saddle points).

## Setup

- 100 trajectories on a spiral (2-dimensional linear ODE)
- Comparison: Levenberg-Marquardt with two complexity metrics as stopping criteria.

## Setup

- 100 trajectories on a spiral (2-dimensional linear ODE)
- Comparison: Levenberg-Marquardt with two complexity metrics as stopping criteria.

$$\begin{tabular}{|c|c|c|c|}\hline Criterion & Best error in $\pmb{X}_*$ \\ \hline $\|\pmb{J}(\pmb{x}_k)^{\mathrm{T}}\pmb{r}(\pmb{x}_k)\| \leq 10^{-3}$ & 49 \\ \hline $\|\underline{\pmb{J}(\pmb{x}_k)^{\mathrm{T}}\pmb{r}(\pmb{x}_k)\|} & \leq 10^{-3}$ or $\|\pmb{r}(\pmb{x}_k)\| \leq 10^{-6}$ & 56. \end{tabular}$$

## Complexity and nonlinear least squares

- A family of complexity metrics and results.
- Derived for Gauss-Newton methods (Levenberg-Marquardt type).
- Works with inexact/stochastic values and derivatives.

E. Bergou, Y. Diouane, V. Kungurstev and C. W. Royer. A stochastic Levenberg-Marquardt method using random models with complexity results. SIAM/ASA JUQ, 2022.

## Complexity and nonlinear least squares

- A family of complexity metrics and results.
- Derived for Gauss-Newton methods (Levenberg-Marquardt type).
- Works with inexact/stochastic values and derivatives.

E. Bergou, Y. Diouane, V. Kungurstev and C. W. Royer. A stochastic Levenberg-Marquardt method using random models with complexity results. SIAM/ASA JUQ, 2022.

## Next

- Gauss-Newton vs Newton steps?
- Application to NeuralODE training.
- Go beyond least squares (cross-entropy loss).

A. Allauzen, I. S. Legheraba and C. W. Royer. Optimization landscape of linear neural ordinary differential equations, in preparation.

## Thank you, and happy birthday Steve!



Harrison Ford	Steve Wright
Part of the Star Wars saga	Part of the IPM saga
Plays a professor/adventurer with a hat and a whip	Is a professor and from Australia

C. W. Royer