

Derivative-free optimisation for least-squares problems

*Joint work with Coralia Cartis (Oxford), Jan Fiala & Benjamin Marteau (NAG Ltd.),
Simon Tett (Edinburgh), Matthias Ehrhardt (Bath)*

Lindon Roberts, ANU (lindon.roberts@anu.edu.au)

Applied Mathematics Seminar, UNSW

16 April 2020

Supported by EPSRC (EP/L015803/1) & NAG Ltd.

1. Introduction to derivative-free optimisation (DFO)
2. DFO for nonlinear least-squares
3. Software implementation
4. Application: parameter tuning of climate models
5. Application: learning image denoising parameters

1. **Introduction to derivative-free optimisation (DFO)**
2. DFO for nonlinear least-squares
3. Software implementation
4. Application: parameter tuning of climate models
5. Application: learning image denoising parameters

Nonlinear Optimisation

Interested in nonlinear, nonconvex optimisation

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}),$$

where objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$.

- Ubiquitous in quantitative disciplines, but very difficult to solve in general
- No information about structure of f , except assumed smoothness
 - Problem constants (e.g. bounds on derivatives) unknown
 - Allow inaccurate evaluation of f , e.g. stochastic noise, iterative process
- Unconstrained, but software allows bounds $a_i \leq x_i \leq b_i$
- Seek local minimiser: $f(\mathbf{x}^*) \leq f(\mathbf{x})$ for all \mathbf{x} close to \mathbf{x}^* (not all $\mathbf{x} \in \mathbb{R}^n$)
 - Actually, seek (approximate) stationary point $\|\nabla f(\mathbf{x}^*)\|_2 \leq \epsilon$
- Gap between 'textbook' algorithms and state-of-the-art performance is large

Basic trust-region method:

- Approximate f near \mathbf{x}_k with **quadratic model**

$$f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{s}$$

- Minimise model (set $\nabla m_k = 0$) gives Newton's method

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k \quad \text{where} \quad [\nabla^2 f(\mathbf{x}_k)] \mathbf{s}_k = -\nabla f(\mathbf{x}_k)$$

but may not converge!

- One way to guarantee convergence: restrict the step size

$$\mathbf{s}_k = \arg \min_{\mathbf{s} \in \mathbb{R}^n} m_k(\mathbf{s}) \quad \text{subject to} \quad \|\mathbf{s}\|_2 \leq \Delta_k$$

\implies 'trust region' subproblem – specialised algorithms exist

Basic iterative method:

1. Given \mathbf{x}_k and $\Delta_k > 0$, evaluate $f(\mathbf{x}_k)$, $\nabla f(\mathbf{x}_k)$, $\nabla^2 f(\mathbf{x}_k)$ and construct model m_k
2. Solve trust region subproblem to get step \mathbf{s}_k
3. Evaluate $f(\mathbf{x}_k + \mathbf{s}_k)$ and determine quality of step

$$\rho_k := \frac{\text{actual decrease}}{\text{predicted decrease}} = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{s}_k)}{m_k(0) - m_k(\mathbf{s}_k)}$$

4. Accept/reject step and update Δ_k :

- If $\rho_k \geq 0.7$, set $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$ and $\Delta_{k+1} = 2\Delta_k$ *[very successful]*
- If $\rho_k \in [0.1, 0.7)$, set $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$ and $\Delta_{k+1} = \Delta_k$ *[successful]*
- If $\rho_k < 0.1$, set $\mathbf{x}_{k+1} = \mathbf{x}_k$ and $\Delta_{k+1} = \Delta_k/2$ *[unsuccessful]*

Standard algorithm with theoretical guarantees (e.g. $\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_k)\|_2 = 0$)

Derivative-Free Optimisation

$$f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{s}$$

- How to calculate derivatives of f in practice?

Derivative-Free Optimisation

$$f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{s}$$

- How to calculate derivatives of f in practice?
 - Write code by hand
 - Finite differences
 - Algorithmic differentiation

Derivative-Free Optimisation

$$f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{s}$$

- How to calculate derivatives of f in practice?
 - Write code by hand
 - Finite differences
 - Algorithmic differentiation
- Difficulties when function evaluation is
 - ‘Black-box’
 - Noisy
 - Computationally expensive

Derivative-Free Optimisation

$$f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{s}$$

- How to calculate derivatives of f in practice?
 - Write code by hand
 - Finite differences
 - Algorithmic differentiation
- Difficulties when function evaluation is
 - ‘Black-box’
 - Noisy
 - Computationally expensive
- Alternative — [derivative-free optimisation \(DFO\)](#)
- Many applications: finance, climate, engineering design, experimental design, ...

Model-Based DFO — Basic Ideas

Many different approaches: [model-based](#), Nelder-Mead, pattern/direct search, genetic algorithms, ...

- Previously,

$$f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{s}$$

- Instead, approximate

$$f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s}) = f(\mathbf{x}_k) + \mathbf{g}_k^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top H_k \mathbf{s}$$

- Find \mathbf{g}_k and H_k without using derivatives

Model-Based DFO — Basic Ideas

Many different approaches: [model-based](#), Nelder-Mead, pattern/direct search, genetic algorithms, ...

- Previously,

$$f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{s}$$

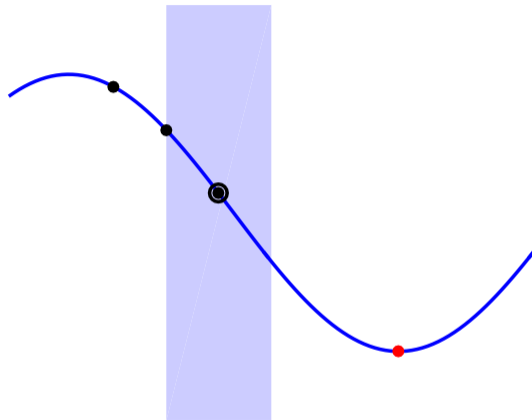
- Instead, approximate

$$f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s}) = f(\mathbf{x}_k) + \mathbf{g}_k^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top H_k \mathbf{s}$$

- Find \mathbf{g}_k and H_k without using derivatives
- How? [Interpolate](#) f over a set of points
- Same trust region approach as before

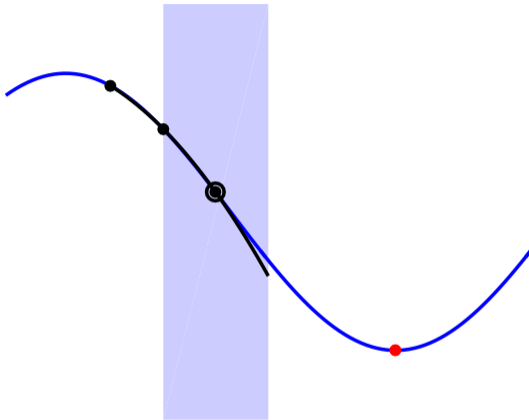
[Conn, Powell, Scheinberg, Toint, Vicente, ...]

Example: Model-Based DFO



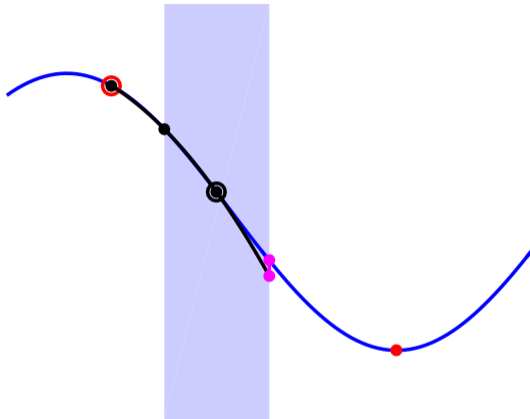
1. Choose interpolation set

Example: Model-Based DFO



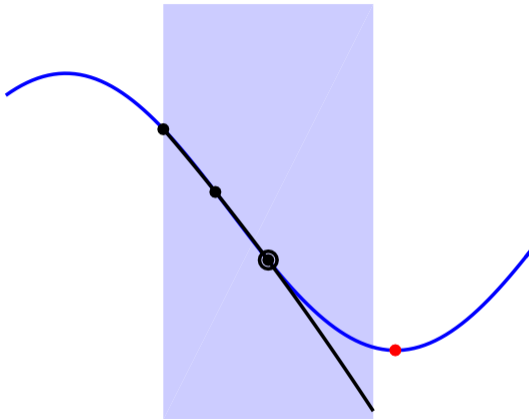
2. Interpolate & minimise...

Example: Model-Based DFO



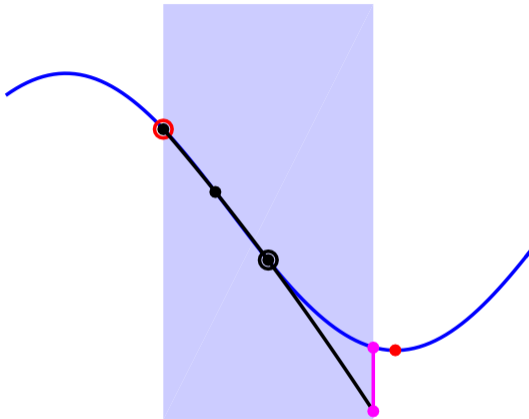
3. Add new point to interpolation set (replace a bad point)

Example: Model-Based DFO



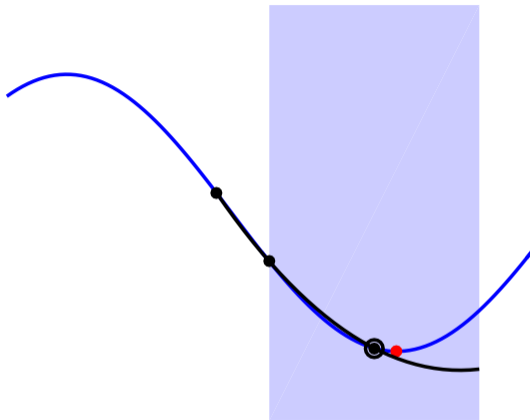
4. Repeat with new interpolation set & model

Example: Model-Based DFO



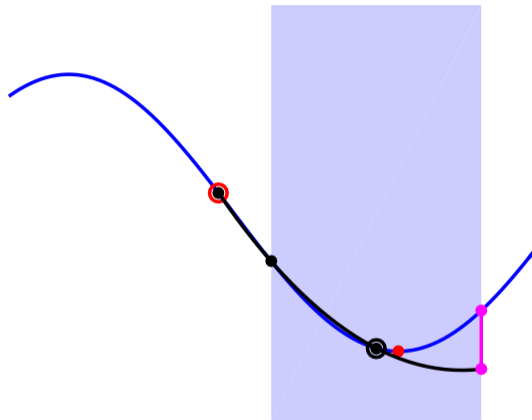
4. Repeat with new interpolation set & model

Example: Model-Based DFO



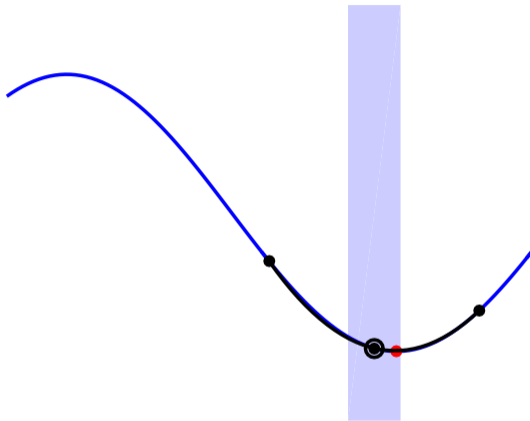
4. Repeat with new interpolation set & model

Example: Model-Based DFO



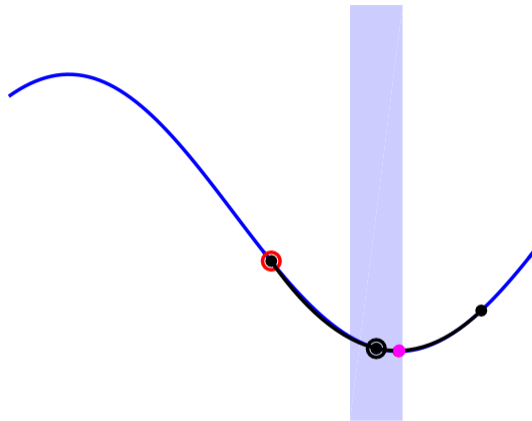
4. Repeat with new interpolation set & model

Example: Model-Based DFO



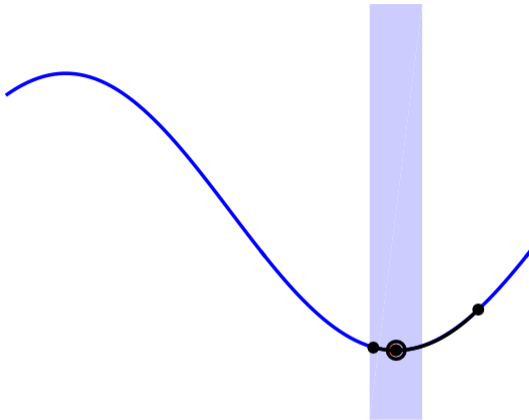
4. Repeat with new interpolation set & model

Example: Model-Based DFO



4. Repeat with new interpolation set & model

Example: Model-Based DFO



4. Repeat with new interpolation set & model

What about theory?

- If geometry of interpolation points is good (in a specific sense), then model has same accuracy order as first-order Taylor series

$$\begin{aligned} |f(\mathbf{x}_k + \mathbf{s}) - m_k(\mathbf{s})| &\leq \kappa_f \Delta_k^2, \\ \|\nabla f(\mathbf{x}_k + \mathbf{s}) - \nabla m_k(\mathbf{s})\|_2 &\leq \kappa_g \Delta_k, \end{aligned}$$

for all $\|\mathbf{s}\|_2 \leq \Delta_k$.

- Need to modify algorithm to fix geometry, when needed
- Get similar convergence results as derivative-based methods
 - Global convergence to stationary points: $\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_k)\|_2 = 0$
 - Worst-case complexity: need at most $\mathcal{O}(\epsilon^{-2})$ iterations to get $\|\nabla f(\mathbf{x}_k)\|_2 \leq \epsilon$

Geometry Requirement

- Interpolation set is $\{\mathbf{y}_0, \dots, \mathbf{y}_p\}$, usually with $\mathbf{y}_0 := \mathbf{x}_k$
- Build model by imposing $f(\mathbf{y}_t) = m_k(\mathbf{y}_t - \mathbf{x}_k)$ for all t
- Lagrange polynomials: $\ell_t(\mathbf{y}_s) = \delta_{s,t}$ for all s, t
- ‘Good’ geometry if all ℓ_t small \Leftrightarrow interpolation problem well-conditioned

Theorem (Conn, Scheinberg & Vicente, 2008)

If $|\ell_t(\mathbf{x}_k + \mathbf{s})| \leq C_1$ for all $\|\mathbf{s}\|_2 \leq \Delta_k$ and all $\|\mathbf{y}_t - \mathbf{x}_k\|_2 \leq C_2\Delta_k$, then the interpolation model is “fully linear” (Taylor-accurate) inside the trust region.

Simple algorithms can check/ensure these conditions (maybe moving some points).

Basic model-based DFO method:

1. Given \mathbf{x}_k and $\Delta_k > 0$, **construct interpolation model** m_k
2. Solve trust region subproblem to get step \mathbf{s}_k
3. Evaluate $f(\mathbf{x}_k + \mathbf{s}_k)$ and determine quality of step

$$\rho_k := \frac{\text{actual decrease}}{\text{predicted decrease}} = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{s}_k)}{m_k(0) - m_k(\mathbf{s}_k)}$$

4. Accept/reject step and update Δ_k :
 - If $\rho_k \geq 0.7$, set $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$ and $\Delta_{k+1} = 2\Delta_k$, **add \mathbf{x}_{k+1} to model** [*very successful*]
 - If $\rho_k \in [0.1, 0.7)$, set $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$ and $\Delta_{k+1} = \Delta_k$, **add \mathbf{x}_{k+1} to model** [*successful*]
 - **If $\rho_k < 0.1$ and model not fully linear, set $\mathbf{x}_{k+1} = \mathbf{x}_k$, $\Delta_{k+1} = \Delta_k$ and make model fully linear** [*model-improving*]
 - **If $\rho_k < 0.1$ and model fully linear, set $\mathbf{x}_{k+1} = \mathbf{x}_k$ and $\Delta_{k+1} = \Delta_k/2$** [*unsuccessful*]

1. Introduction to derivative-free optimisation (DFO)
2. **DFO for nonlinear least-squares**
3. Software implementation
4. Application: parameter tuning of climate models
5. Application: learning image denoising parameters

DFO for Least-Squares — Basic Framework

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|_2^2, \quad \mathbf{r}(\mathbf{x}) \in \mathbb{R}^m$$

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|_2^2, \quad \mathbf{r}(\mathbf{x}) \in \mathbb{R}^m$$

Classical Gauss-Newton

- Linearise \mathbf{r} at \mathbf{x}_k using Jacobian

$$\mathbf{r}(\mathbf{x}_k + \mathbf{s}) \approx \mathbf{m}_k(\mathbf{s}) = \mathbf{r}(\mathbf{x}_k) + \mathbf{J}(\mathbf{x}_k)\mathbf{s}$$

- Approximation: $\nabla^2 f(\mathbf{x}_k) =$
 $\mathbf{J}(\mathbf{x}_k)^\top \mathbf{J}(\mathbf{x}_k) + \sum_{i=1}^m r_i(\mathbf{x}_k) \nabla^2 r_i(\mathbf{x}_k)$

Derivative-Free Gauss-Newton

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|_2^2, \quad \mathbf{r}(\mathbf{x}) \in \mathbb{R}^m$$

Classical Gauss-Newton

- Linearise \mathbf{r} at \mathbf{x}_k using Jacobian

$$\mathbf{r}(\mathbf{x}_k + \mathbf{s}) \approx \mathbf{m}_k(\mathbf{s}) = \mathbf{r}(\mathbf{x}_k) + \mathbf{J}(\mathbf{x}_k)\mathbf{s}$$

- Approximation: $\nabla^2 f(\mathbf{x}_k) =$
 $\mathbf{J}(\mathbf{x}_k)^\top \mathbf{J}(\mathbf{x}_k) + \sum_{i=1}^m r_i(\mathbf{x}_k) \nabla^2 r_i(\mathbf{x}_k)$

Derivative-Free Gauss-Newton

- Jacobian not available, use

$$\mathbf{m}_k(\mathbf{s}) = \mathbf{r}(\mathbf{x}_k) + \mathbf{J}_k \mathbf{s}$$

- Find \mathbf{J}_k by interpolation

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|_2^2, \quad \mathbf{r}(\mathbf{x}) \in \mathbb{R}^m$$

Classical Gauss-Newton

- Linearise \mathbf{r} at \mathbf{x}_k using Jacobian

$$\mathbf{r}(\mathbf{x}_k + \mathbf{s}) \approx \mathbf{m}_k(\mathbf{s}) = \mathbf{r}(\mathbf{x}_k) + \mathbf{J}(\mathbf{x}_k)\mathbf{s}$$

- Approximation: $\nabla^2 f(\mathbf{x}_k) =$
 $\mathbf{J}(\mathbf{x}_k)^\top \mathbf{J}(\mathbf{x}_k) + \sum_{i=1}^m r_i(\mathbf{x}_k) \nabla^2 r_i(\mathbf{x}_k)$

Derivative-Free Gauss-Newton

- Jacobian not available, use

$$\mathbf{m}_k(\mathbf{s}) = \mathbf{r}(\mathbf{x}_k) + \mathbf{J}_k \mathbf{s}$$

- Find \mathbf{J}_k by interpolation

In both cases, solve trust region subproblem with **simplified quadratic** model

$$f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s}) = \frac{1}{2} \|\mathbf{m}_k(\mathbf{s})\|_2^2$$

Previous works use **quadratic** models for $\mathbf{r}(\mathbf{x})$ [Zhang, Conn & Scheinberg (2010), Wild (2017)]

Advantages of linear models

- Match global convergence guarantees
- Fewer evaluations of $\mathbf{r}(\mathbf{x})$ to build first model
- Lower linear algebra cost ($\approx 7\times$ speedup) and improved scalability
- Explicit connection between geometry and linear algebra
 - When adding interpolation point \mathbf{x}_{k+1} , delete point \mathbf{y}_t with large $|\ell_t(\mathbf{x}_{k+1})|$
 - Gives rank-1 update of interpolation matrix A , hence

$$A_{\text{new}}^{-1} = A_{\text{old}}^{-1} + \frac{1}{\sigma_t} \mathbf{u}_t \mathbf{v}_t^{\top} \quad \text{where } \sigma_t = \ell_t(\mathbf{x}_{k+1})$$

1. Introduction to derivative-free optimisation (DFO)
2. DFO for nonlinear least-squares
3. **Software implementation**
4. Application: parameter tuning of climate models
5. Application: learning image denoising parameters

DFO-LS (Derivative-Free Optimisation for Least-Squares)

Open-source Python package (NLLS with bounds)

- Github: `numerical algorithms group/dfols`

Key Features:

- **Flexible model construction**
 \implies *enables both reduced initialisation cost and regression models*
- **Robust to noisy objectives** using multiple restarts
 \implies *effective alternative to sample averaging, regression models*
- **Reduced initialisation cost** for expensive objectives
 \implies *progress from 2 evaluations, if desired*

Flexible Model Construction

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|_2^2, \quad \mathbf{r}(\mathbf{x}) \in \mathbb{R}^m$$

- Have $p + 1$ interpolation points $\{\mathbf{y}_0 = \mathbf{x}_k, \mathbf{y}_1, \dots, \mathbf{y}_p\}$
- Find model $\mathbf{m}_k(\mathbf{s}) = \mathbf{r}_k + J_k \mathbf{s}$ by solving

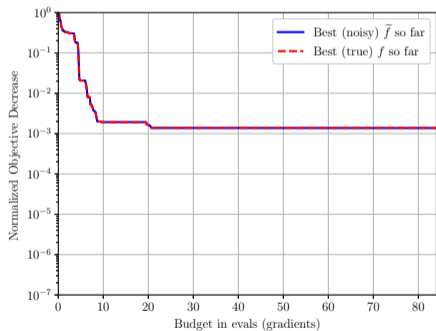
$$\min_{\mathbf{r}_k, J_k} \sum_{t=0}^p \|\mathbf{m}_k(\mathbf{y}_t - \mathbf{x}_k) - \mathbf{r}(\mathbf{y}_t)\|_2^2$$

\Rightarrow one $(p + 1) \times (n + 1)$ system, different RHS for each residual r_i , $i = 1, \dots, m$

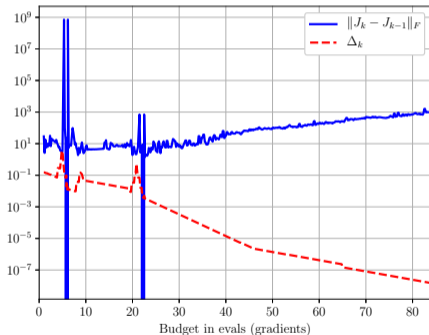
- Works for any $p \geq 1$
 - Unique interpolant if $p = n$ (usual case)
 - Regression model if $p > n$ (e.g. noisy objective)
 - Select minimal-norm solution if $p < n$ (used for reduced initialisation cost)

Noisy Problems — Example of Stagnation

- TR radius $\Delta_k \rightarrow 0$, so interpolation points eventually get close together
- Objective values all within noise level \implies interpolated model only captures noise
- This is one of the main use cases for DFO!



Normalised Objective Decrease



Convergence Details

Common approaches:

- Sample averaging [Deng & Ferris (2006), Chen, Menickelly & Scheinberg (2016)]
- Regression models [Conn, Scheinberg & Vicente (2009), Billups, Larson & Graf (2013)]
- Both available in DFO-LS if desired (c.f. flexible model construction)

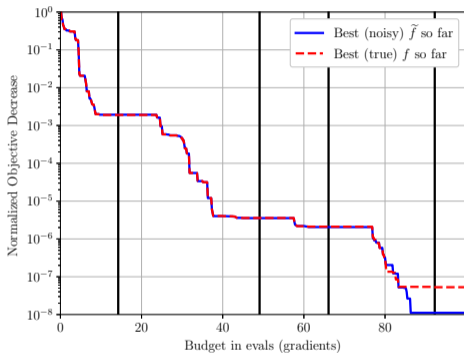
Common approaches:

- Sample averaging [Deng & Ferris (2006), Chen, Menickelly & Scheinberg (2016)]
- Regression models [Conn, Scheinberg & Vicente (2009), Billups, Larson & Graf (2013)]
- Both available in DFO-LS if desired (c.f. flexible model construction)

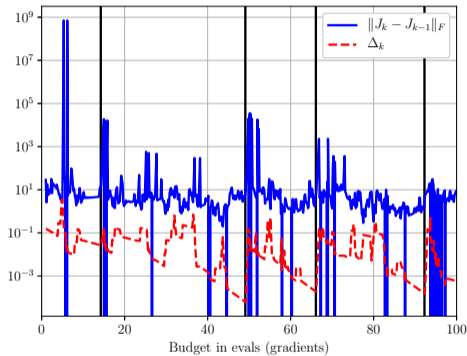
Alternative strategy: multiple restarts

- When $\Delta_k \leq \Delta_{min}$, reset $\Delta_{k+1} = \Delta_0$ initial TR radius
- Update interpolation set for new trust region:
 - Move \mathbf{x}_k plus $N \approx 2$ points closest to \mathbf{x}_k to geometry-improving points in $B(\mathbf{x}_k, \Delta_{k+1})$
- *Auto-detection*: call restart if Δ_k consistently decreasing and $\|J_k - J_{k-1}\|_F$ consistently increasing for several iterations

Multiple Restarts Strategy — Example



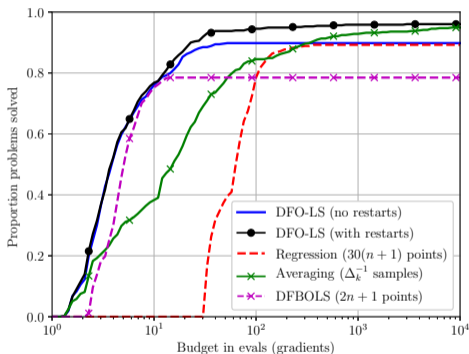
Normalised Objective Decrease



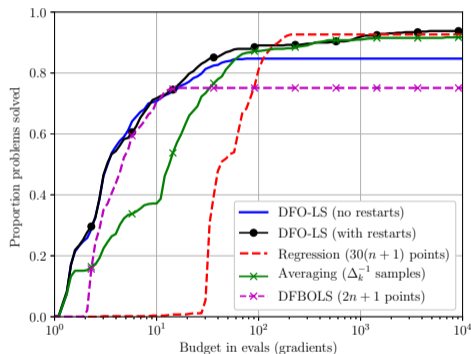
Convergence Details

DFO-LS — Comparison of Noise Robustness Strategies

Data profiles (using 53 test problems from [Moré & Wild, 2009], $\tau = \max(10^{-5}, \hat{\tau})$)



Mult. Gaussian noise ($\pm 1\%$)



Add. Gaussian noise (± 0.01)

% test problems solved with a given # objective evaluations; higher values are better

Reduced Initialisation Cost for Expensive Objectives

- Start by evaluating objective at \mathbf{x}_0 and p random orthogonal directions $\mathbf{y}_t = \mathbf{x}_0 + \Delta_0 \mathbf{q}_t$ ($t = 1, \dots, p$) — usually $p = n$
- **Expensive objective?** Want to see progress with very few evaluations

Reduced Initialisation Cost for Expensive Objectives

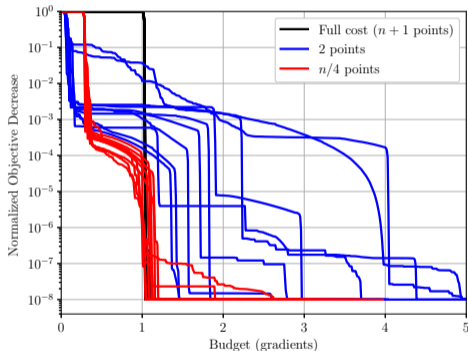
- Start by evaluating objective at \mathbf{x}_0 and p random orthogonal directions $\mathbf{y}_t = \mathbf{x}_0 + \Delta_0 \mathbf{q}_t$ ($t = 1, \dots, p$) — usually $p = n$
- **Expensive objective?** Want to see progress with very few evaluations
- For $p < n$ directions, use interpolating model with minimal norm
- **Problem:** J_k not full rank, so $\mathbf{s}_k \in \text{span}\{\mathbf{y}_1 - \mathbf{x}_k, \dots, \mathbf{y}_p - \mathbf{x}_k\}$
⇒ can never search outside the initial subspace of directions

Reduced Initialisation Cost for Expensive Objectives

- Start by evaluating objective at \mathbf{x}_0 and p random orthogonal directions $\mathbf{y}_t = \mathbf{x}_0 + \Delta_0 \mathbf{q}_t$ ($t = 1, \dots, p$) — usually $p = n$
- **Expensive objective?** Want to see progress with very few evaluations
- For $p < n$ directions, use interpolating model with minimal norm
- **Problem:** J_k not full rank, so $\mathbf{s}_k \in \text{span}\{\mathbf{y}_1 - \mathbf{x}_k, \dots, \mathbf{y}_p - \mathbf{x}_k\}$
⇒ can never search outside the initial subspace of directions
- **Solution:** Artificially perturb J_k to make it full rank
 - Floor singular values at $\sigma_p > 0$
- Sometimes this will give descent, but always expands the search space

Reduced Initialisation Cost — Example

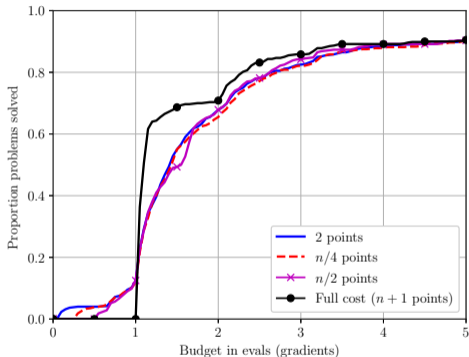
Can make reasonable progress with $< n + 1$ evaluations, but usually better to wait (if possible)



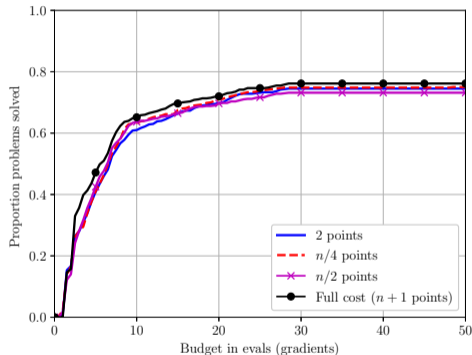
Objective decrease for BROWNALE ($n = 100$)

Reduced Initialisation Cost — Performance

Data profiles (60 test problems with $n \approx 100$)



Short budget, $\tau = 10^{-1}$



Long budget, $\tau = 10^{-5}$

% test problems solved with a given # objective evaluations; higher values are better

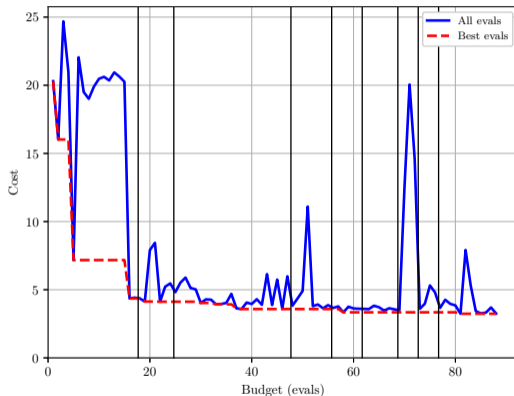
1. Introduction to derivative-free optimisation (DFO)
2. DFO for nonlinear least-squares
3. Software implementation
4. **Application: parameter tuning of climate models**
5. Application: learning image denoising parameters

Climate Parameter Tuning

- Aim: tuning models of global atmospheric physics
- Fit to observations (e.g. average temperature, humidity, radiation)
- Difficulty: simulations are expensive (multi-year global climate simulation) and noisy (underlying physics is chaotic)
- Standard approach for tuning climate models is manual:
 - Generate different sets of parameters
 - Evaluate fit to observations
 - Select parameters with best fit (and possibly perturb these to generate new sets)
- Alternative: apply DFO-LS with multiple restarts

Climate Parameter Tuning — Results

Example Results (HadAM3, 14 parameters, budget 90 evaluations)



Cost vs. # objective evaluations (cost ≤ 5 considered broadly in line with observations)

Climate Parameter Tuning — Results

- Start DFO-LS from 5 different starting locations
- Find 5 different parameter combinations, all in line with observations (and genuinely different local minima)
- Outperforms other solvers: approximate finite differencing, surrogate modelling

Climate research question

Q: How are these local minima different from a climatology perspective?

A: Not very! *Uncertainty in climate predictions largely driven by modelling choices (processes included & parametrisations) not parameter tuning.*

1. Introduction to derivative-free optimisation (DFO)
2. DFO for nonlinear least-squares
3. Software implementation
4. Application: parameter tuning of climate models
5. **Application: learning image denoising parameters**

Many image processing problems can be posed in the form

$$\min_{\mathbf{x}} \mathcal{D}(A\mathbf{x}, \mathbf{y}) + \mathcal{R}(\mathbf{x}),$$

where \mathcal{D} measures data fidelity ($A\mathbf{x} \approx \mathbf{y}$) and \mathcal{R} is regulariser; e.g. denoising

$$\min_{\mathbf{x}} \underbrace{\frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2}_{\mathcal{D}(\mathbf{x}, \mathbf{y})} + \alpha \underbrace{\sum_j \sqrt{\|\nabla x_j\|_2^2 + \epsilon^2}}_{\approx \text{TV}(\mathbf{x})} + \frac{\eta}{2} \|\mathbf{x}\|_2^2$$

This problem is **smooth and strongly convex**, can be solved effectively with iterative methods (gradient descent, NAG, FISTA, etc.).

Bilevel Learning

- Unclear how to choose parameters $\theta := [\alpha, \epsilon, \eta]^\top$
- One option: learn parameters from example problems $\{(\mathbf{x}_i, \mathbf{y}_i)\}$:

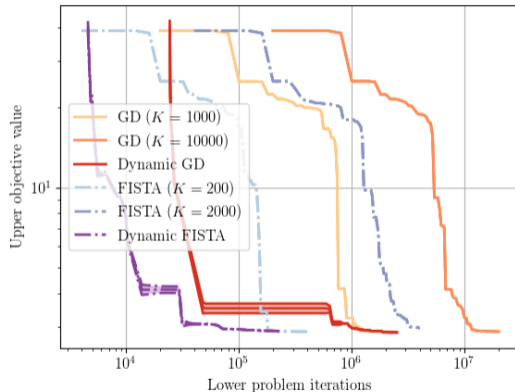
$$\min_{\theta} \sum_i \|\hat{\mathbf{x}}_i(\theta) - \mathbf{x}_i\|_2^2$$

s.t. $\hat{\mathbf{x}}_i$ solves denoising problem with θ

- Bilevel optimisation problem, requires computing $\partial_{\theta} \hat{\mathbf{x}}_i(\theta)$
 - Requires very high accuracy solves of denoising problem
 - Don't know in advance what accuracy is required (educated guess)
- **Alternative:** modify DFO-LS to allow **dynamic accuracy** on objective evaluations (i.e. ask for $\hat{\mathbf{x}}_i(\theta)$ correct to within some error δ_x)

Bilevel Learning

Example results (gradient descent & FISTA as lower-level solvers):





More efficient learning, without requiring heuristics for lower-level accuracy




Conclusions




- DFO methods suitable when objective is expensive and/or noisy
- DFO equivalent of Gauss-Newton gives an effective algorithm for least-squares
- Reduced initialisation cost if desired — can start progressing after 2 evaluations
- Effective for tuning global climate models and bilevel learning


Future work:

- Scalability: dimensionality reduction, sparsity, inexact interpolation solves, etc.
- Local convergence rates
- General objective and constrained problems

-  Stephen C. Billups, Jeffrey W. Larson, and Peter Graf.
Derivative-free optimization of expensive functions with computational error using weighted regression.
SIAM Journal on Optimization, 23(1):27–53, 2013.
-  Coralia Cartis, Jan Fiala, Benjamin Marteau, and Lindon Roberts.
Improving the flexibility and robustness of model-based derivative-free optimization solvers.
ACM Transactions on Mathematical Software, 45(3):32:1–32:41, 2019.

-  Coralia Cartis and Lindon Roberts.
A derivative-free Gauss-Newton method.
Mathematical Programming Computation, 2019.
-  Ruobing Chen, Matt Menickelly, and Katya Scheinberg.
Stochastic optimization using a trust-region method and random models.
Mathematical Programming, 169(2):447–487, 2018.
-  Andrew R. Conn, Katya Scheinberg, and Luís N. Vicente.
Geometry of interpolation sets in derivative free optimization.
Mathematical Programming, 111(1-2):141–172, 2007.

-  Andrew R. Conn, Katya Scheinberg, and Luís N. Vicente.
Introduction to Derivative-Free Optimization, volume 8 of MPS-SIAM Series on Optimization.
MPS/SIAM, Philadelphia, 2009.
-  Matthias J. Ehrhardt and Lindon Roberts.
Inexact derivative free optimization for bi-level learning.
in preparation, 2020.
-  Simon F. B. Tett, Jonathan M. Gregory, Nicolas Freychet, Coralia Cartis, Michael J. Mineter, and Lindon Roberts.
Does model calibration reduce uncertainty in climate projections?
submitted, 2020.

-  Stefan M. Wild.
POUNDERS in TAO: Solving derivative-free nonlinear least-squares problems with POUNDERS.
In T. Terlaky, M. F. Anjos, and S. Ahmed, editors, *Advances and Trends in Optimization with Engineering Applications*, volume 24 of *MOS-SIAM Book Series on Optimization*, pages 529–539. MOS/SIAM, Philadelphia, 2017.
-  Hongchao Zhang, Andrew R. Conn, and Katya Scheinberg.
A derivative-free algorithm for least-squares minimization.
SIAM Journal on Optimization, 20(6):3555–3576, 2010.