

# Large-scale derivative-free optimization using random subspace methods

*Joint work with Coralia Cartis (Oxford) & Clément Royer (Paris-Dauphine PSL)*

---

Lindon Roberts, University of Sydney ([lindon.roberts@sydney.edu.au](mailto:lindon.roberts@sydney.edu.au))

Simons Collaboration on Hidden Symmetries and Fusion Energy

1 September 2023

## Further Reading

This talk is based on:

- C. Cartis & L. Roberts, Scalable subspace methods for derivative-free nonlinear least-squares optimization, *Math. Prog.*, 2023.
- L. Roberts & C. W. Royer, Direct search based on probabilistic descent in reduced spaces, *SIAM J. Optim.*, to appear.

Our software packages are:

- DFBN for nonlinear least-squares:  
<https://github.com/numericalalgorithmsgroup/dfbgn>
- directsearch for general problems:  
<https://github.com/lindonroberts/directsearch>

1. **Introduction to derivative-free optimization (DFO)**
2. Subspace DFO methods
3. Numerical results

# Nonlinear Optimization

Interested in **unconstrained nonlinear optimization**

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

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

- $f$  is possibly nonconvex and/or 'black-box'
  - In practice, allow inaccurate evaluations of  $f$ , e.g. noise, outcome of iterative process
- Seek **local minimizer** (actually, approximate stationary point:  $\|\nabla f(\mathbf{x})\|_2 \leq \epsilon$ )

Lots of high-quality algorithms available:

- Linesearch,  $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k H_k^{-1} \nabla f(\mathbf{x}_k)$  (e.g. GD, Newton, BFGS)
- **Trust-region methods** (adapt well to derivative-free setting)
- Others: cubic regularization, nonlinear CG, ...

## Basic trust-region method

- Approximate  $f$  near  $\mathbf{x}_k$  with a **local quadratic (Taylor) model**

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

- Get step by minimizing model in a neighborhood

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

- Accept/reject step and adjust  $\Delta_k$  based on quality of new point  $f(\mathbf{x}_k + \mathbf{s}_k)$

$$\mathbf{x}_{k+1} = \begin{cases} \mathbf{x}_k + \mathbf{s}_k, & \text{if sufficient decrease,} & \longleftarrow (\text{maybe increase } \Delta_k) \\ \mathbf{x}_k, & \text{otherwise.} & \longleftarrow (\text{decrease } \Delta_k) \end{cases}$$

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

[Conn, Gould & Toint, 2000]

# Derivative-Free Optimization

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

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

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

# Derivative-Free Optimization

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

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

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

# Derivative-Free Optimization

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

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

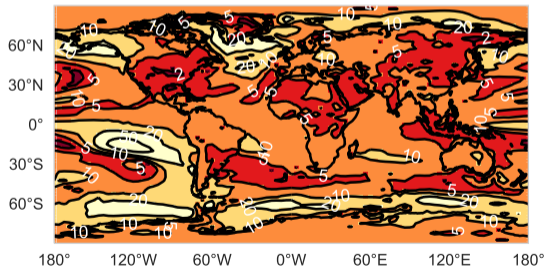
- How to calculate derivatives of  $f$  in practice?
  - Write code by hand
  - Finite differences
  - Algorithmic differentiation/backpropagation
- Difficulties when function evaluation is
  - Black-box
  - Noisy
  - Computationally expensive
- Alternative — derivative-free optimization (DFO)



## Application 1: Climate Modelling

[Tett et al., 2022]

- Parameter calibration for global climate models
- One model run = simulate global climate for 5 years (expensive!)
- Very complicated, chaotic physics (black-box & noisy!)



## Application 2: Adversarial Example Generation

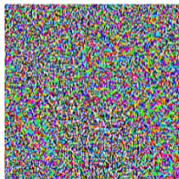
[Alzantot et al., 2019]

- Find perturbations of neural network inputs which are misclassified
- Neural network structure assumed to be unknown (black-box!)
- Want to test very few examples ( $\approx$  expensive!)



“panda”  
57.7% confidence

+ .007 ×



=



“gibbon”  
99.3 % confidence

*Image from [Goodfellow et al., 2015]*

## DFO Method 1: Model-Based DFO

- Using trust-region framework, build a model

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

and find  $\mathbf{g}_k$  and  $\mathbf{H}_k$  without using derivatives

## DFO Method 1: Model-Based DFO

- Using trust-region framework, build a model

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

and find  $\mathbf{g}_k$  and  $\mathbf{H}_k$  without using derivatives

- How? **Interpolate  $f$  over a set of points** — find  $\mathbf{g}_k, \mathbf{H}_k$  such that

$$m_k(\mathbf{y} - \mathbf{x}_k) = f(\mathbf{y}), \quad \forall \mathbf{y} \in \mathcal{Y}$$

## DFO Method 1: Model-Based DFO

- Using trust-region framework, build a model

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

and find  $\mathbf{g}_k$  and  $\mathbf{H}_k$  without using derivatives

- How? **Interpolate  $f$  over a set of points** — find  $\mathbf{g}_k, \mathbf{H}_k$  such that

$$m_k(\mathbf{y} - \mathbf{x}_k) = f(\mathbf{y}), \quad \forall \mathbf{y} \in \mathcal{Y}$$

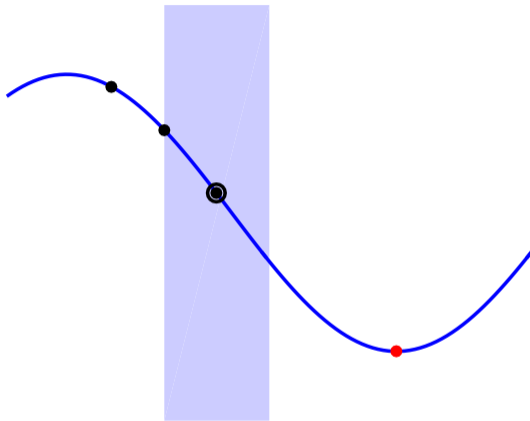
For convergence, need  $m_k$  to be **fully linear**:

$$|f(\mathbf{x}_k + \mathbf{s}) - m_k(\mathbf{s})| \leq \mathcal{O}(\Delta_k^2) \quad \text{and} \quad \|\nabla f(\mathbf{x}_k + \mathbf{s}) - \nabla m_k(\mathbf{s})\|_2 \leq \mathcal{O}(\Delta_k)$$

Achievable if points in  $\mathcal{Y}$  are well-spaced (in a specific sense).

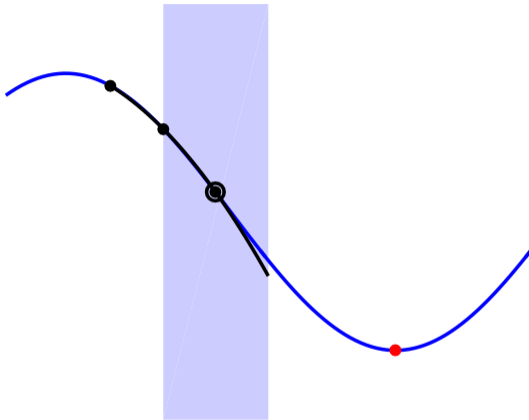
[Powell, 2003; Conn, Scheinberg & Vicente, 2009]

## Example: Model-Based DFO



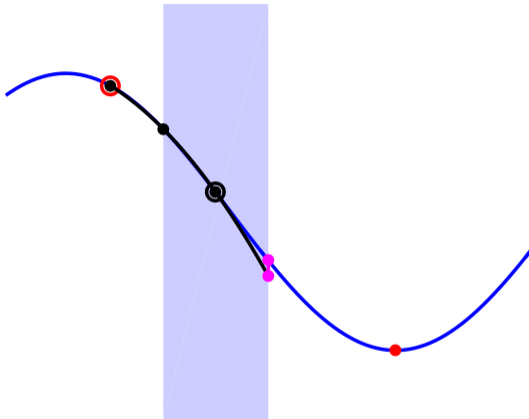
### 1. Choose interpolation set

## Example: Model-Based DFO



### 2. Interpolate & minimize...

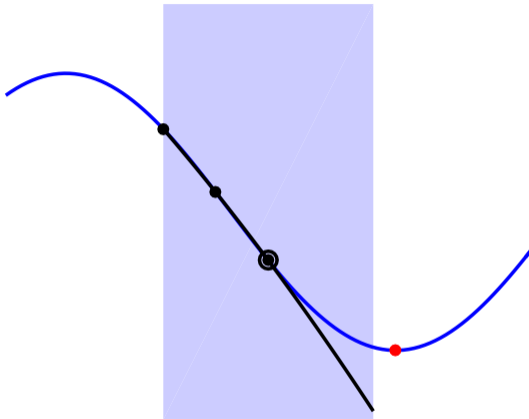
## Example: Model-Based DFO



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

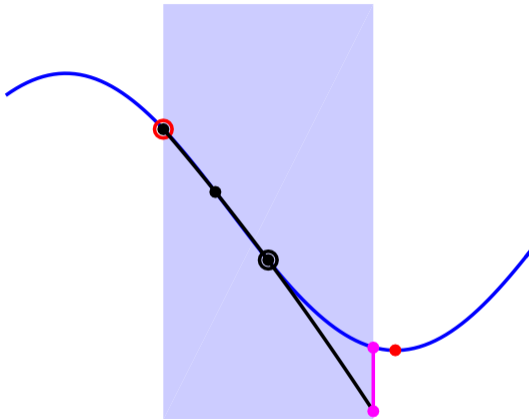


## Example: Model-Based DFO



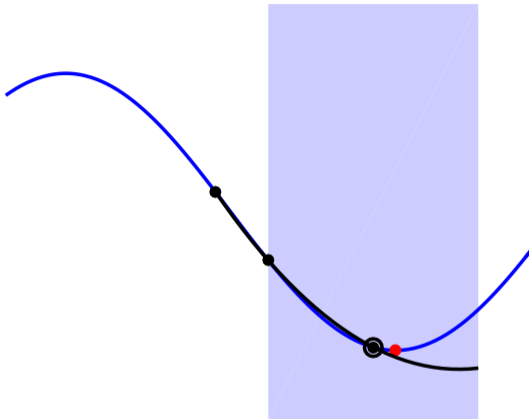
**4. Repeat with new interpolation set & model**

## Example: Model-Based DFO



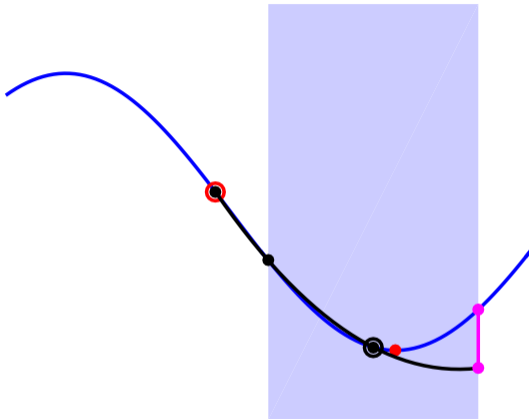
**4. Repeat with new interpolation set & model**

## Example: Model-Based DFO



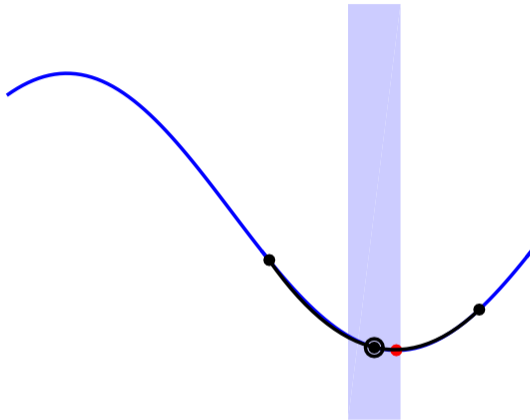
**4. Repeat with new interpolation set & model**

## Example: Model-Based DFO



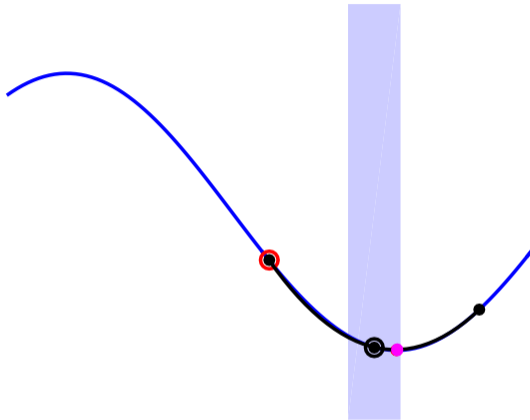
**4. Repeat with new interpolation set & model**

## Example: Model-Based DFO



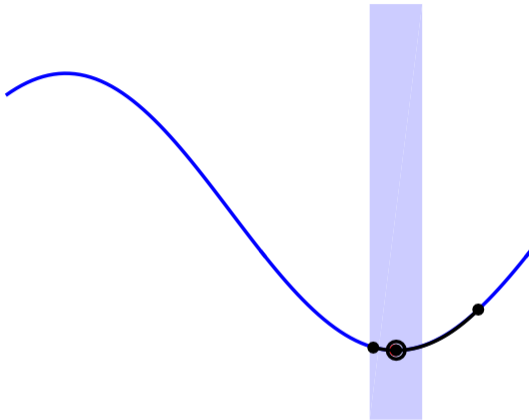
**4. Repeat with new interpolation set & model**

## Example: Model-Based DFO



**4. Repeat with new interpolation set & model**

## Example: Model-Based DFO



**4. Repeat with new interpolation set & model**

## DFO Method 2: Direct Search

- Given  $\mathbf{x}_k$  and  $\Delta_k$ , choose a set  $\mathcal{D}_k \subset \mathbb{R}^n$  of  $m$  vectors
- If there exists  $\mathbf{d}_k \in \mathcal{D}_k$  with  $f(\mathbf{x}_k + \Delta_k \mathbf{d}_k) < f(\mathbf{x}_k) - \frac{1}{2} \Delta_k^2 \|\mathbf{d}_k\|_2^2$ :
  - Set  $\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta_k \mathbf{d}_k$  and increase  $\Delta_k$
- Otherwise, set  $\mathbf{x}_{k+1} = \mathbf{x}_k$  and decrease  $\Delta_k$



## DFO Method 2: Direct Search

- Given  $\mathbf{x}_k$  and  $\Delta_k$ , choose a set  $\mathcal{D}_k \subset \mathbb{R}^n$  of  $m$  vectors
- If there exists  $\mathbf{d}_k \in \mathcal{D}_k$  with  $f(\mathbf{x}_k + \Delta_k \mathbf{d}_k) < f(\mathbf{x}_k) - \frac{1}{2} \Delta_k^2 \|\mathbf{d}_k\|_2^2$ :
  - Set  $\mathbf{x}_{k+1} = \mathbf{x}_k + \Delta_k \mathbf{d}_k$  and increase  $\Delta_k$
- Otherwise, set  $\mathbf{x}_{k+1} = \mathbf{x}_k$  and decrease  $\Delta_k$

For convergence, need  $\mathcal{D}_k$  to be  $\kappa$ -descent:

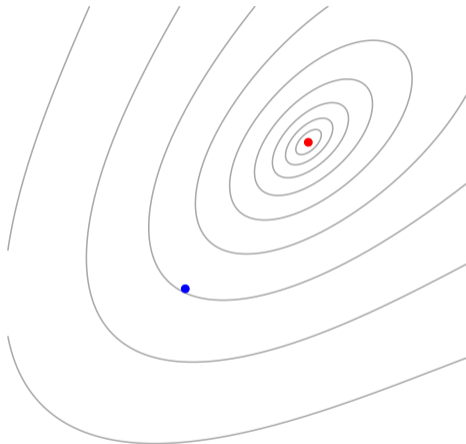
$$\max_{\mathbf{d} \in \mathcal{D}_k} \frac{-\mathbf{d}^T \nabla f(\mathbf{x}_k)}{\|\mathbf{d}\|_2 \cdot \|\nabla f(\mathbf{x}_k)\|_2} \geq \kappa \in (0, 1]$$

i.e. there is a vector  $\mathbf{d}$  making an acute angle with  $-\nabla f(\mathbf{x}_k)$  (descent direction).

Examples:  $\{\pm \mathbf{e}_1, \dots, \pm \mathbf{e}_n\}$  with  $\kappa = 1/\sqrt{n}$  or  $\{\mathbf{e}_1, \dots, \mathbf{e}_n, -\mathbf{e}\}$  with  $\kappa \sim 1/n$ .

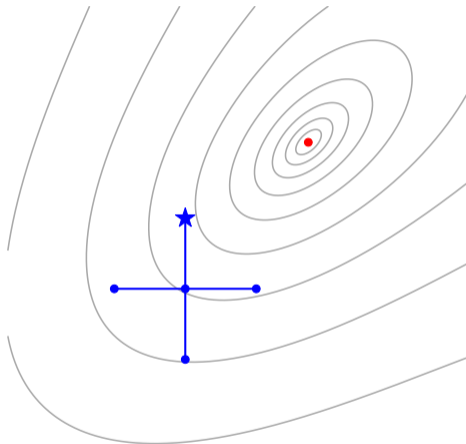
[Kolda, Lewis & Torczon, 2003; Conn, Scheinberg & Vicente, 2009]

## Example: Direct Search



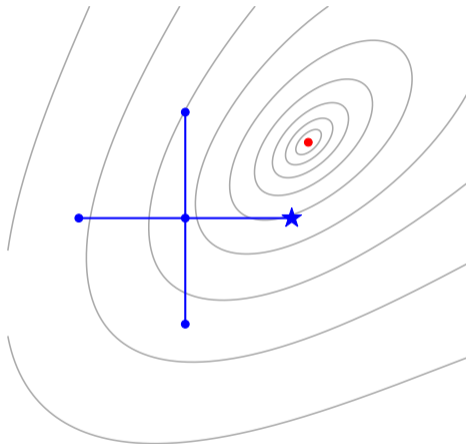
Modified from [Kolda, Lewis & Torczon, 2003]

## Example: Direct Search



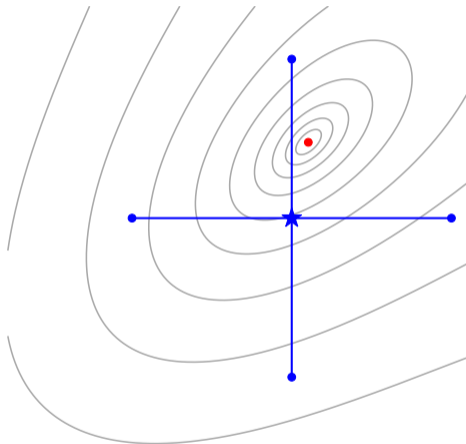
Modified from [Kolda, Lewis & Torczon, 2003]

## Example: Direct Search



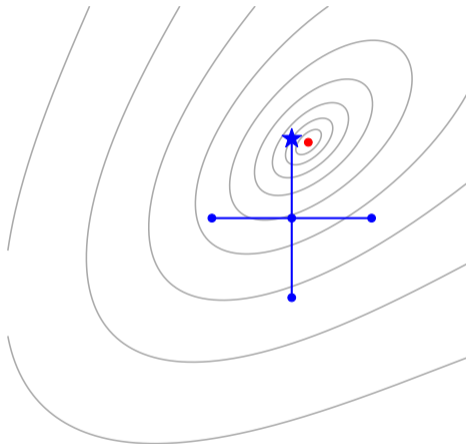
Modified from [Kolda, Lewis & Torczon, 2003]

## Example: Direct Search



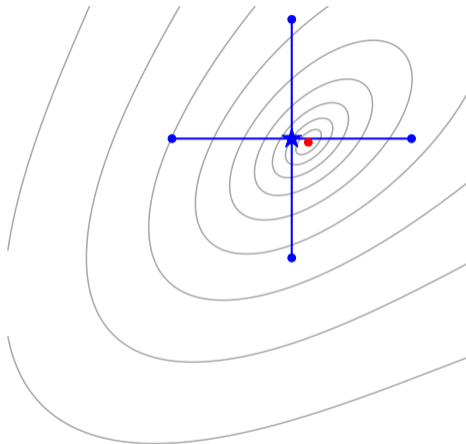
Modified from [Kolda, Lewis & Torczon, 2003]

## Example: Direct Search



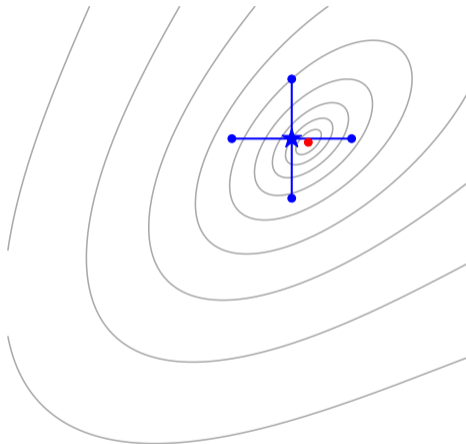
Modified from [Kolda, Lewis & Torczon, 2003]

## Example: Direct Search



Modified from [Kolda, Lewis & Torczon, 2003]

## Example: Direct Search



Modified from [Kolda, Lewis & Torczon, 2003]



Analyze methods using **worst-case complexity**: how long before  $\|\nabla f(\mathbf{x}_k)\|_2 \leq \epsilon$ ?

Metric	Deriv-based	Model-based	Direct search
Iterations	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(n^2\epsilon^{-2})$	$\mathcal{O}(n\epsilon^{-2})$
Evaluations	$\approx \mathcal{O}(n\epsilon^{-2})$	$\mathcal{O}(n^3\epsilon^{-2})$	$\mathcal{O}(n^2\epsilon^{-2})$

[Cartis, Gould & Toint, 2010; Garmanjani, Júdice & Vicente, 2016; Vicente, 2013]

- Same  $\epsilon$  dependency as derivative-based, but **scales badly with problem dimension  $n$**
- Model-based DFO also has substantial linear algebra work for interpolation and geometry management: at least  $\mathcal{O}(n^3)$  flops per iteration

## Challenge

How can DFO methods be made scalable?

Analyze methods using **worst-case complexity**: how long before  $\|\nabla f(\mathbf{x}_k)\|_2 \leq \epsilon$ ?

Metric	Deriv-based	Model-based	Direct search
Iterations	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(n^2\epsilon^{-2})$	$\mathcal{O}(n\epsilon^{-2})$
Evaluations	$\approx \mathcal{O}(n\epsilon^{-2})$	<del><math>\mathcal{O}(n^3\epsilon^{-2})</math></del> $\mathcal{O}(n^2\epsilon^{-2})$	<del><math>\mathcal{O}(n^2\epsilon^{-2})</math></del> $\mathcal{O}(n\epsilon^{-2})$

[Cartis, Gould & Toint, 2010; Garmanjani, Júdice & Vicente, 2016; Vicente, 2013]

- Same  $\epsilon$  dependency as derivative-based, but ~~scales badly with problem dimension  $n$~~
- Model-based DFO also has substantial linear algebra work for interpolation and geometry management: at least  ~~$\mathcal{O}(n^3)$~~   $\mathcal{O}(n)$  flops per iteration

## Challenge

How can DFO methods be made scalable?

1. Introduction to derivative-free optimization (DFO)
2. **Subspace DFO methods**
3. Numerical results

## Challenge

How can DFO methods be made scalable?

- Exploit known problem structure [Porcelli & Toint, 2020; Bandeira et al., 2012]
- Randomized finite differencing ('gradient sampling') [Nesterov & Spokoiny, 2017]

Applications for scalable DFO methods include:

- Machine learning [Salimans et al., 2017; Ughi et al., 2020]
- Image analysis [Ehrhardt & R., 2021]
- Proxy for global optimization methods [Cartis, R. & Sheridan-Methven, 2021]

## Challenge

How can DFO methods be made scalable?

Randomization is a promising approach:

- Make model fully linear with probability  $< 1$  [Gratton et al., 2017]
- Make search directions  $\kappa$ -descent with probability  $< 1$  [Gratton et al., 2015]

## Challenge

How can DFO methods be made scalable?

Randomization is a promising approach:

- Make model fully linear with probability  $< 1$  [Gratton et al., 2017]
- Make search directions  $\kappa$ -descent with probability  $< 1$  [Gratton et al., 2015]

**Problem:** Improves complexity for direct search, but not for model-based!

**Why?** Direct search formulation effectively allows dimensionality reduction (sample  $\ll n$  directions).

## Goal

Use dimensionality reduction techniques suitable for both DFO classes.

### Lemma (Johnson-Lindenstrauss, 1984)

Suppose  $X$  is a set of  $N$  points in  $\mathbb{R}^d$  and  $\epsilon \in (0, 1)$ . Let  $A \in \mathbb{R}^{p \times d}$  be a matrix with i.i.d.  $N(0, p^{-2})$  entries and  $p \sim \log(N)/\epsilon$ . Then with high probability,

$$(1 - \epsilon)\|x - y\|_2 \leq \|Ax - Ay\|_2 \leq (1 + \epsilon)\|x - y\|_2, \quad \forall x, y \in X.$$

### Lemma (Johnson-Lindenstrauss, 1984)

Suppose  $X$  is a set of  $N$  points in  $\mathbb{R}^d$  and  $\epsilon \in (0, 1)$ . Let  $A \in \mathbb{R}^{p \times d}$  be a matrix with i.i.d.  $N(0, p^{-2})$  entries and  $p \sim \log(N)/\epsilon$ . Then with high probability,

$$(1 - \epsilon)\|x - y\|_2 \leq \|Ax - Ay\|_2 \leq (1 + \epsilon)\|x - y\|_2, \quad \forall x, y \in X.$$

- Random projections approximately preserve distances (& inner products, norms, ...)
- Reduced dimension  $p$  depends only on # of points  $N$ , **not the ambient dimension  $d$ !**
- Other random constructions satisfy J-L Lemma (Haar subsampling, hashing, ...)



# Subspace DFO

We use a subspace method: only search in **low-dimensional subspaces** of  $\mathbb{R}^n$

- Related to coordinate descent methods [Wright, 2015; Patrascu & Necoara, 2015]
- Some implementations exist, but no theory [Gross & Parks, 2020; Neumaier et al., 2011]
- Build on recent derivative-based analysis [Cartis, Fowkes & Shao, 2020]

# Subspace DFO

We use a subspace method: only search in **low-dimensional subspaces** of  $\mathbb{R}^n$

- Related to coordinate descent methods [Wright, 2015; Patrascu & Necoara, 2015]
- Some implementations exist, but no theory [Gross & Parks, 2020; Neumaier et al., 2011]
- Build on recent derivative-based analysis [Cartis, Fowkes & Shao, 2020]

## Subspace DFO framework:

- Generate subspace of dimension  $p \ll n$  given by  $\text{col}(P_k)$  for random  $P_k \in \mathbb{R}^{n \times p}$
- Model-based: build a low-dimensional model  $\hat{m}_k(\hat{\mathbf{s}})$  which is fully linear for  $\hat{f}(\hat{\mathbf{s}}) := f(\mathbf{x}_k + P_k \hat{\mathbf{s}}) : \mathbb{R}^p \rightarrow \mathbb{R}$
- Direct search: choose  $\mathcal{D}_k \subset \mathbb{R}^p$  which is  $\kappa$ -descent for  $P_k^T \nabla f(\mathbf{x}_k) \in \mathbb{R}^p$

Fewer interpolation/sample points needed, cheap linear algebra (everything in  $\mathbb{R}^p$ )

## Subspace DFO — Subspace Quality

**Choice of subspace:** we need to make sure we search in ‘good’ subspaces (where there is potential to decrease  $f$  sufficiently).

The subspace at iteration  $k$  is **well-aligned** if

$$\|P_k^T \nabla f(\mathbf{x}_k)\|_2 \geq \alpha \|\nabla f(\mathbf{x}_k)\|_2, \quad \text{for some } \alpha > 0.$$

i.e. if there is still work to do, then we know this by only inspecting  $f$  in the subspace.

## Subspace DFO — Subspace Quality

**Choice of subspace:** we need to make sure we search in ‘good’ subspaces (where there is potential to decrease  $f$  sufficiently).

The subspace at iteration  $k$  is **well-aligned** if

$$\|P_k^T \nabla f(\mathbf{x}_k)\|_2 \geq \alpha \|\nabla f(\mathbf{x}_k)\|_2, \quad \text{for some } \alpha > 0.$$

i.e. if there is still work to do, then we know this by only inspecting  $f$  in the subspace.

### Key Assumption

The subspace  $P_k$  is well-aligned with probability  $1 - \delta$ .

Using J-L lemma, choose  $p \sim (1 - \alpha)^{-2} |\log \delta| = \mathcal{O}(1)$  independent of  $n$ .

*Note: if randomly select  $p$  coordinates (block coordinate descent), need  $p \sim \alpha n$ .*

### Theorem (Cartis & R., 2023; R. & Royer, 2023)

*If  $f$  is sufficiently smooth and bounded below and  $\epsilon$  sufficiently small, then*

$$\mathbb{P} [K_\epsilon \leq C(p, \alpha, \delta)\epsilon^{-2}] \geq 1 - e^{-c(p, \alpha, \delta)\epsilon^{-2}},$$

*where  $K_\epsilon$  is the first iteration with  $\|\nabla f(\mathbf{x}_k)\|_2 \leq \epsilon$ .*

- Implies  $\mathbb{E} [K_\epsilon] = \mathcal{O}(\epsilon^{-2})$  and almost-sure convergence
- $\mathcal{O}(p)$  evaluations per iteration, so same bounds for evaluation complexity

## Standard methods:

Metric	Deriv-based	Model-based	Direct search
Iterations	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(n^2\epsilon^{-2})$	$\mathcal{O}(n\epsilon^{-2})$
Evaluations	$\approx \mathcal{O}(n\epsilon^{-2})$	$\mathcal{O}(n^3\epsilon^{-2})$	$\mathcal{O}(n^2\epsilon^{-2})$

Model-based DFO has  $\mathcal{O}(n^3)$  linear algebra work per iteration.

## Using random subspaces:

Metric	Deriv-based	Model-based	Direct search
Iterations	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(n^2\epsilon^{-2})$	$\mathcal{O}(n\epsilon^{-2})$
Evaluations	$\approx \mathcal{O}(n\epsilon^{-2})$	$\mathcal{O}(n^2\epsilon^{-2})$	$\mathcal{O}(n\epsilon^{-2})$

Model-based DFO has  $\mathcal{O}(n)$  linear algebra work per iteration.

1. Introduction to derivative-free optimization (DFO)
2. Subspace DFO methods
3. **Numerical results**

# Software Packages

Open-source Python packages available on Github

## Model-Based

DFBGN for nonlinear least-squares (`numericalalgorithmsgroup/dfbgn`)

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|_2^2 = \frac{1}{2} \sum_{i=1}^m r_i(\mathbf{x})^2$$

Subspace method with several heuristics to improve performance

## Direct Search

`directsearch` (`lindonroberts/directsearch`)

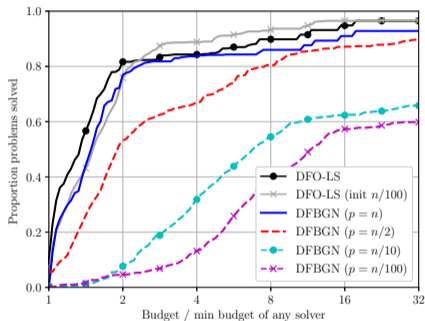
Many varieties of direct search methods (classical, random, subspaces) with multiple  $\mathcal{D}_k$  generation methods.



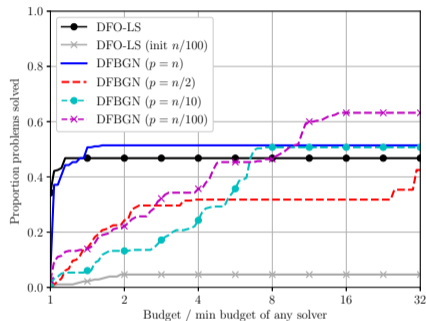
# Numerical Results — DFBGN

DFBGN vs. DFO-LS (low accuracy  $\tau = 10^{-1}$ )

[% problems solved vs. # evals]



Medium-scale problems,  $n \approx 100$



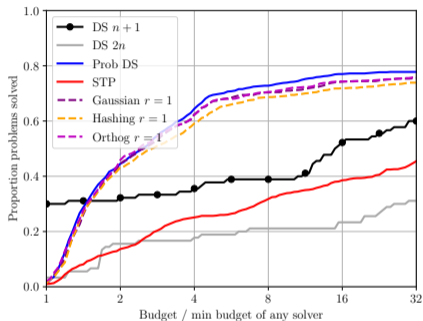
Large problems  $n \approx 1000$ , 12hr timeout

DFBGN is more suitable for low accuracy solutions, performance improves with larger  $p$  (except for timeouts!)

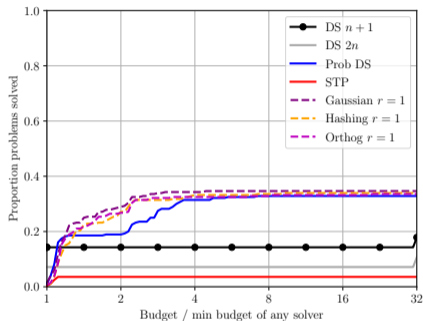
# Numerical Results — Direct Search

Direct search comparisons (low accuracy  $\tau = 10^{-1}$ )

[% problems solved vs. # evals]



Medium-scale problems,  $n \approx 100$

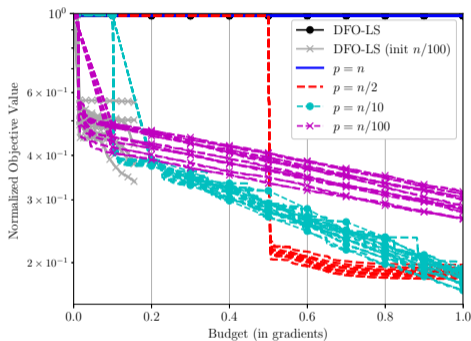


Large problems  $n \approx 1000$

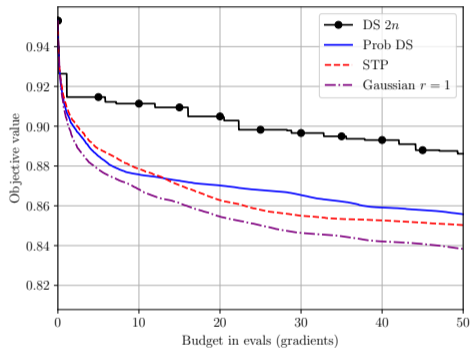
Subspace methods match randomized methods and outperform classical methods, performance best with small  $p$

# Numerical Results — low budget

Subspace methods progress after  $p \ll n$  evaluations (important when  $n$  large)



**DFBGN**



**directsearch**

*(normalized objective reduction vs. # evaluations, 12hr timeout)*

## Conclusions

- Scalability of model-based DFO is currently limited (in theory & practice)
- Randomized projections are effective for dimensionality reduction
- New algorithms reduce linear algebra cost and iteration complexity
- Practical implementations available

## Future Work

- Second-order complexity analysis
- Efficient implementation of subspace quadratic models (model-based)
- Problems with constraints
- Comparison of different choices of  $p$ :
  - New work ( $\sim 3$  weeks ago!) studying this [Hare, R. & Royer, 2023]

- M. ALZANTOT, Y. SHARMA, S. CHAKRABORTY, H. ZHANG, C.-J. HSIEH, AND M. B. SRIVASTAVA, *GenAttack: Practical black-box attacks with gradient-free optimization*, in Proceedings of the Genetic and Evolutionary Computation Conference, Prague, Czech Republic, 2019, ACM, pp. 1111–1119.
- A. S. BANDEIRA, K. SCHEINBERG, AND L. N. VICENTE, *Computation of sparse low degree interpolating polynomials and their application to derivative-free optimization*, Mathematical Programming, 134 (2012), pp. 223–257.
- C. CARTIS, J. FOWKES, AND Z. SHAO, *A randomised subspace Gauss-Newton method for nonlinear least-squares*, in Workshop on “Beyond first-order methods in ML systems” at the 37th International Conference on Machine Learning, Vienna, Austria, 2020.
- C. CARTIS, N. I. M. GOULD, AND P. L. TOINT, *On the complexity of steepest descent, Newton’s and regularized Newton’s methods for nonconvex unconstrained optimization problems*, SIAM Journal on Optimization, 20 (2010), pp. 2833–2852.
- C. CARTIS AND L. ROBERTS, *Scalable subspace methods for derivative-free nonlinear least-squares optimization*, Mathematical Programming, 199 (2023), pp. 461—524.

- C. CARTIS, L. ROBERTS, AND O. SHERIDAN-METHVEN, *Escaping local minima with local derivative-free methods: A numerical investigation*, Optimization, (2021).
- A. R. CONN, N. I. M. GOULD, AND P. L. TOINT, *Trust-Region Methods*, vol. 1 of MPS-SIAM Series on Optimization, MPS/SIAM, Philadelphia, 2000.
- A. R. CONN, K. SCHEINBERG, AND L. N. VICENTE, *Introduction to Derivative-Free Optimization*, vol. 8 of MPS-SIAM Series on Optimization, MPS/SIAM, Philadelphia, 2009.
- M. J. EHRHARDT AND L. ROBERTS, *Inexact derivative-free optimization for bilevel learning*, Journal of Mathematical Imaging and Vision, 63 (2020), pp. 580–600.
- R. GARMANJANI, D. JÚDICE, AND L. N. VICENTE, *Trust-region methods without using derivatives: Worst case complexity and the nonsmooth case*, SIAM Journal on Optimization, 26 (2016), pp. 1987–2011.
- I. J. GOODFELLOW, J. SHLENS, AND C. SZEGEDY, *Explaining and harnessing adversarial examples*, in 3rd International Conference on Learning Representations ICLR, San Diego, 2015.
- S. GRATTON, C. W. ROYER, L. N. VICENTE, AND Z. ZHANG, *Direct search based on probabilistic descent*, SIAM Journal on Optimization, 25 (2015), pp. 1515–1541.

- S. GRATTON, C. W. ROYER, L. N. VICENTE, AND Z. ZHANG, *Complexity and global rates of trust-region methods based on probabilistic models*, IMA Journal of Numerical Analysis, 38 (2017), pp. 1579–1597.
- J. C. GROSS AND G. T. PARKS, *Optimization by moving ridge functions: Derivative-free optimization for computationally intensive functions*, arXiv preprint arXiv:2007.04893, (2020).
- W. HARE, L. ROBERTS, AND C. W. ROYER, *Expected decrease for derivative-free algorithms using random subspaces*, arXiv preprint arXiv:2308.04734, (2023).
- T. G. KOLDA, R. M. LEWIS, AND V. TORCZON, *Optimization by direct search: New perspectives on some classical and modern methods*, SIAM Review, 45 (2003), pp. 385–482.
- Y. NESTEROV AND V. SPOKOINY, *Random gradient-free minimization of convex functions*, Foundations of Computational Mathematics, 17 (2017), pp. 527–566.
- A. NEUMAIER, H. FENDL, H. SCHILLY, AND T. LEITNER, *VXQR: Derivative-free unconstrained optimization based on QR factorizations*, Soft Computing, 15 (2011), pp. 2287–2298.
- A. PATRASCU AND I. NECOARA, *Efficient random coordinate descent algorithms for large-scale structured nonconvex optimization*, Journal of Global Optimization, 61 (2015), pp. 19–46.

- M. PORCELLI AND P. L. TOINT, *Global and local information in structured derivative free optimization with BFO*, arXiv preprint arXiv:2001.04801, (2020).
- M. J. D. POWELL, *On trust region methods for unconstrained minimization without derivatives*, Mathematical Programming, 97 (2003), pp. 605–623.
- L. ROBERTS AND C. W. ROYER, *Direct search based on probabilistic descent in reduced spaces*, SIAM Journal on Optimization, to appear (2023).
- T. SALIMANS, J. HO, X. CHEN, S. SIDOR, AND I. SUTSKEVER, *Evolution strategies as a scalable alternative to reinforcement learning*, arXiv preprint arXiv:1703.03864, (2017).
- S. F. B. TETT, J. M. GREGORY, N. FREYCHET, C. CARTIS, M. J. MINETER, AND L. ROBERTS, *Does model calibration reduce uncertainty in climate projections?*, Journal of Climate, 35 (2022), pp. 2585–2602.
- G. UGHI, V. ABROL, AND J. TANNER, *An empirical study of derivative-free-optimization algorithms for targeted black-box attacks in deep neural networks*, arXiv preprint arXiv:2012.01901, (2020).
- L. N. VICENTE, *Worst case complexity of direct search*, EURO Journal on Computational Optimization, 1 (2013), pp. 143–153.



S. J. WRIGHT, *Coordinate descent algorithms*, *Mathematical Programming*, 151 (2015), pp. 3–34.