

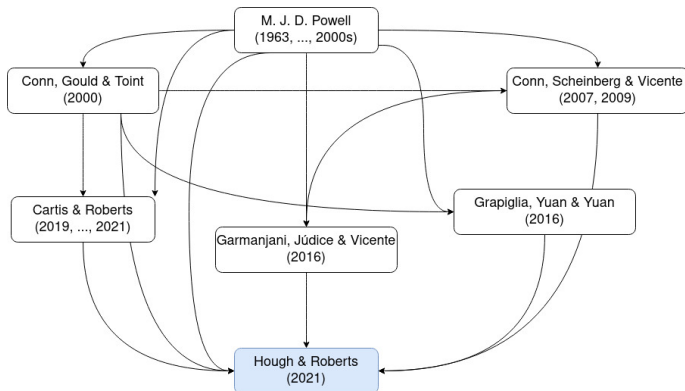
Model-Based Derivative-Free Methods for Constrained Optimization

Joint work with Lindon Roberts (ANU)

Matthew Hough, University of Queensland \Rightarrow University of Waterloo
mhough@uwaterloo.ca

November 1st, 2021

Background



Outline

1. Introduction to DFO trust-region methods
2. Handling constraints
3. Application to composite minimization
4. Numerical results

The Problem

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

- ▶ $f : \mathbb{R}^n \rightarrow \mathbb{R}$ continuously differentiable and possibly nonconvex
- ▶ Assume we cannot evaluate $\nabla f(\mathbf{x})$
 - ▶ Black-box
 - ▶ Noisy
 - ▶ Computationally expensive
- ▶ Applications: climate modelling, experimental design, machine learning, etc
- ▶ Seeking a local minimizer (approx. stationary point: $\|\nabla f(\mathbf{x}^*)\| \leq \epsilon$)

Model-Based DFO

- ▶ Classic approach:

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

- ▶ Instead, approximate:

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

- ▶ Find \mathbf{g}_k and H_k by interpolating f over a set of points

Model-Based DFO: Algorithm

(assuming our interpolation model is a good approx.)

1. Build local interpolation model:

$$f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s})$$

2. Minimize the model within the trust-region Δ_k to get the step

$$\mathbf{s}_k = \arg \min_{\mathbf{s} \in \mathbb{R}^d} m_k(\mathbf{s}) \quad \text{s.t.} \quad \|\mathbf{s}\|_2 \leq \Delta_k$$

3. Evaluate $f(\mathbf{x}_k + \mathbf{s}_k)$, check sufficient decrease, select \mathbf{x}_{k+1} and Δ_{k+1}
4. Update interpolation set with the new point $\mathbf{x}_k + \mathbf{s}_k$

Model-Based DFO: Interpolation Geometry

We may not get sufficient decrease if...

1. Δ_k is too large
2. m_k is not a good approximation to f (bad geometry)

Problems!

- ▶ How to ensure good geometry?
- ▶ How do we define good geometry?

Good geometry \implies accurate model \implies convergence

Model-Based DFO: Interpolation Geometry

- ▶ Need interpolation set $\{\mathbf{y}_0, \dots, \mathbf{y}_n\}$ to be "well-poised" in $B(\mathbf{y}_0, \Delta)$
- ▶ **Λ -poised** if all $\mathbf{y}_t \in B(\mathbf{y}_0, \Delta)$ and exists $\Lambda \geq 1$ s.t.

$$\max |\ell_t(\mathbf{y})| \leq \Lambda, \quad \forall \mathbf{y} \in B(\mathbf{y}_0, \Delta)$$

- ▶ $\ell_t(\mathbf{y}_s) = \delta_{s,t}$ for all s, t
- ▶ Points are "well-spaced"

[Conn, Scheinberg & Vicente, 2009]

Model-based DFO: Interpolation Geometry

- ▶ Λ -poisedness \implies fully linear model:

- ▶ $|f(\mathbf{x}_k + \mathbf{s}) - m(\mathbf{s})| \leq \kappa_{ef} \Delta_k^2$

(κ_{ef}, κ_{eg} depend on Λ)

- ▶ $\|\nabla f(\mathbf{x}_k + \mathbf{s}) - \nabla m(\mathbf{s})\| \leq \kappa_{eg} \Delta_k$

for all $\mathbf{y} \in B(\mathbf{y}_0, \Delta_k)$, $\|\mathbf{s}\| \leq \Delta_k$

- ▶ Fully linear model \implies convergence

- ▶ Two important algorithms:

1. Checks $\{\mathbf{y}_0, \dots, \mathbf{y}_n\}$ is Λ -poised
2. Makes $\{\mathbf{y}_0, \dots, \mathbf{y}_n\}$ Λ -poised if it is not already

[Conn, Scheinberg & Vicente, 2009]

The Constrained Problem

$$\min_{\mathbf{x} \in \mathcal{C}} f(\mathbf{x})$$

- ▶ $f : \mathbb{R}^n \rightarrow \mathbb{R}$ continuously differentiable and possibly nonconvex
- ▶ Assume we cannot evaluate $\nabla f(\mathbf{x})$
- ▶ $\mathcal{C} \subseteq \mathbb{R}^d$ has nonempty interior, closed, and convex
 - ▶ Cannot evaluate f outside of \mathcal{C}
 - ▶ Only accessible via projection, $P_{\mathcal{C}} : \mathbb{R}^d \rightarrow \mathcal{C}$

Constrained DFO: Algorithm

1. Build local interpolation model from feasible points:

$$f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s})$$

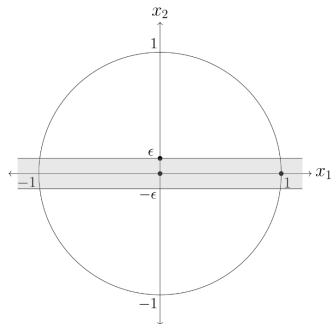
2. Minimize the model within $B(\mathbf{y}_0, \Delta_k) \cap \mathcal{C}$ to get the step

$$\mathbf{s}_k = \arg \min_{\mathbf{s} \in B(\mathbf{y}_0, \Delta_k) \cap \mathcal{C}} m_k(\mathbf{s})$$

3. Evaluate $f(\mathbf{x}_k + \mathbf{s}_k)$, check sufficient decrease, select \mathbf{x}_{k+1} and Δ_{k+1}
4. Update interpolation set with the new point $\mathbf{x}_k + \mathbf{s}_k$

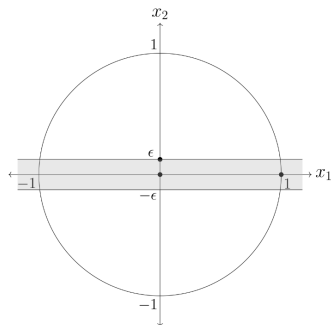
Constrained DFO: Geometry

- ▶ $C = \{(x_1, x_2) : |x_2| \leq \epsilon\} \subseteq \mathbb{R}^2$
- ▶ $Y = \{(0, 0), (1, 0), (0, \epsilon)\} \subseteq B(\mathbf{0}, 1)$
- ▶ In $B(\mathbf{0}, 1)$, points are Λ -poised with $\Lambda \sim \epsilon^{-1}$
 \implies large κ_{ef}, κ_{eg}



Constrained DFO: Geometry

- ▶ $\mathcal{C} = \{(x_1, x_2) : |x_2| \leq \epsilon\} \subseteq \mathbb{R}^2$
- ▶ $Y = \{(0, 0), (1, 0), (0, \epsilon)\} \subseteq B(\mathbf{0}, 1)$
- ▶ In $B(\mathbf{0}, 1)$, points are Λ -poised with $\Lambda \sim \epsilon^{-1}$
 \implies large κ_{ef}, κ_{eg}



- ▶ Λ -poised if all $\mathbf{y}_t \in B(\mathbf{y}_0, \Delta) \cap \mathcal{C}$ and exists $\Lambda \geq 1$ s.t.

$$\max |\ell_t(\mathbf{y})| \leq \Lambda, \quad \forall \mathbf{y} \in B(\mathbf{y}_0, \Delta) \cap \mathcal{C}$$

- ▶ Now we have $\Lambda \leq 3$ independent of $\epsilon \implies$ improved error bounds

Constrained DFO: Geometry

- ▶ Λ -poisedness \implies fully linear model in $B(\mathbf{x}_k, \Delta_k)$:

$$\max_{\substack{\mathbf{x}_k + \mathbf{s} \in \mathcal{C} \\ \|\mathbf{s}\| \leq \Delta_k}} |f(\mathbf{x}_k + \mathbf{s}) - m_k(\mathbf{s})| \leq \kappa_{ef} \Delta_k^2$$

$$\max_{\substack{\mathbf{x}_k + \mathbf{s} \in \mathcal{C} \\ \|\mathbf{s}\| \leq 1}} |(\nabla f(\mathbf{x}_k) - \mathbf{g}_k)^T \mathbf{s}| \leq \kappa_{eg} \Delta_k$$

- ▶ Slightly weaker:
 - ▶ $\nabla m(\mathbf{y}) \approx \nabla f(\mathbf{y})$ only at $\mathbf{y} = \mathbf{x}_k$
 - ▶ Only care about points in \mathcal{C}
- ▶ Still have important algorithms
 1. Check points are Λ -poised
 2. Make points Λ -poised if not

Constrained DFO: Measuring Progress

$$\pi^f(\mathbf{x}) := \left| \min_{\substack{\mathbf{x}+\mathbf{s} \in \mathcal{C} \\ \|\mathbf{s}\| \leq 1}} \nabla f(\mathbf{x})^T \mathbf{s} \right| \quad \Longrightarrow \quad \pi^m(\mathbf{x}) := \left| \min_{\substack{\mathbf{x}+\mathbf{s} \in \mathcal{C} \\ \|\mathbf{s}\| \leq 1}} \mathbf{g}_k^T \mathbf{s} \right|$$

- ▶ For $\mathcal{C} = \mathbb{R}^d$, $\pi^f(\mathbf{x}_k) = \|\nabla f(\mathbf{x}_k)\|$, and $\pi^g(\mathbf{x}_k) = \|\mathbf{g}_k\|$
- ▶ fully linear $\Longrightarrow |\pi^f(\mathbf{x}_k) - \pi^m(\mathbf{x}_k)| \leq \kappa_{eg} \Delta_k$

[Conn, Gould & Toint, 2000]

Constrained DFO: Measuring Progress

$$\pi^f(\mathbf{x}) := \left| \min_{\substack{\mathbf{x}+\mathbf{s} \in \mathcal{C} \\ \|\mathbf{s}\| \leq 1}} \nabla f(\mathbf{x})^T \mathbf{s} \right| \quad \Longrightarrow \quad \pi^m(\mathbf{x}) := \left| \min_{\substack{\mathbf{x}+\mathbf{s} \in \mathcal{C} \\ \|\mathbf{s}\| \leq 1}} \mathbf{g}_k^T \mathbf{s} \right|$$

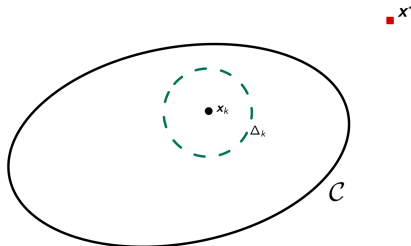
- ▶ For $\mathcal{C} = \mathbb{R}^d$, $\pi^f(\mathbf{x}_k) = \|\nabla f(\mathbf{x}_k)\|$, and $\pi^g(\mathbf{x}_k) = \|\mathbf{g}_k\|$
- ▶ fully linear $\Longrightarrow |\pi^f(\mathbf{x}_k) - \pi^m(\mathbf{x}_k)| \leq \kappa_{eg} \Delta_k$

Solution to $\pi^f(\mathbf{x})$ is given by $\mathbf{s}^* := p(t, \mathbf{x}) - \mathbf{x}$

- ▶ where $p(t, \mathbf{x}) = P_{\mathcal{C}}(\mathbf{x} - t\nabla f(\mathbf{x}))$, $t \geq 0$,
- ▶ and $\|p(t, \mathbf{x}) - \mathbf{x}\| = 1$

[Conn, Gould & Toint, 2000]

Constrained DFO: Measuring Progress

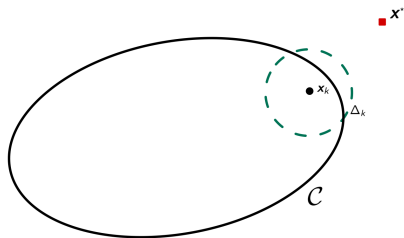


$$\blacktriangleright \mathbf{s}^* = P_C(\mathbf{x} - t\nabla f(\mathbf{x})) - \mathbf{x} = -t\nabla f(\mathbf{x})$$

$$\blacktriangleright 1 = \|\mathbf{s}^*\| = t\|\nabla f(\mathbf{x})\| \implies t = \frac{1}{\|\nabla f(\mathbf{x})\|} \implies \mathbf{s}^* = \frac{-\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|}$$

$$\blacktriangleright \implies \pi^f(\mathbf{x}) = \|\nabla f(\mathbf{x})\|$$

Constrained DFO: Measuring Progress



- ▶ $P_C(\mathbf{x} - t\nabla f(\mathbf{x})) \neq \mathbf{x} - t\nabla f(\mathbf{x})$
- ▶ $\mathbf{s}^* = p(t, \mathbf{x}) - \mathbf{x}$ gets smaller near the boundary
- ▶ $\implies \pi^f(\mathbf{x}) \rightarrow 0$ as approach constraints in direction of \mathbf{x}^*

Constrained DFO: Convergence Theory

1. Ensure we always have m_k fully linear (by ensuring good geometry)
 2. Ensure $\pi_k^m \sim \Delta_k$
 3. When $\pi^m(\mathbf{x}_k) \rightarrow 0$, we are also getting $\pi^f(\mathbf{x}_k) \rightarrow 0$
 4. Standard convergence results follow
- Worst-case complexity: at most $\mathcal{O}(\epsilon^{-2})$ iterations to have $\pi_k^m \leq \epsilon$

Application to composite minimization

$$f(\mathbf{x}) = F(\mathbf{r}(\mathbf{x}))$$

- ▶ where $\mathbf{r} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a black-box function
- ▶ Derivatives of $\mathbf{r}(\mathbf{x})$ are unavailable
- ▶ Classic example is $F(\mathbf{r}) = \frac{1}{2}\|\mathbf{r}\|^2$

Application to composite minimization

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

- Typically linearize \mathbf{r} at \mathbf{x}_k using the Jacobian:

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

- But in DFO, Jacobian is not available:

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

- Find \mathbf{J}_k by interpolation

End up with a local quadratic model

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

Implementation

Open-source Python implementation: *DFO-LS*

- ▶ Github: [numericalalgorithmsgroup/dfols](https://github.com/numericalalgorithmsgroup/dfols)
-

- ▶ Replace gradient-descent step with projected gradient-descent (PGD)
- ▶ Dykstra's algorithm for projecting onto \mathcal{C}
- ▶ New point becomes

$$\mathbf{x}_{k+1} = P_{\mathcal{Q}}(\mathbf{x}_k - t\mathbf{g}_k)$$

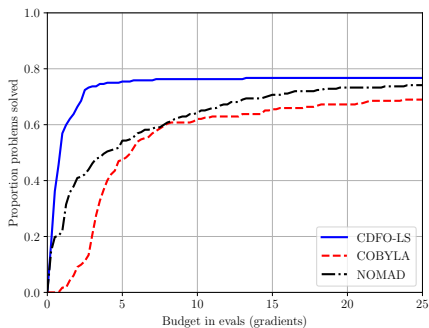
$$\mathcal{Q} := \mathcal{C} \cap B(\mathbf{x}_k, \Delta_k)$$

[Beck, 2017]

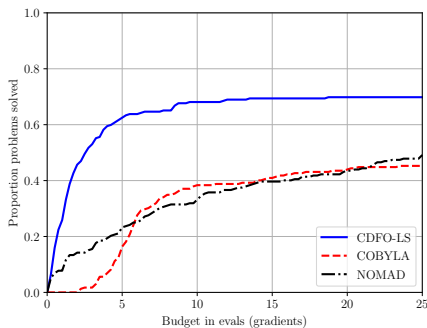
Numerical results

58 test problems with ball, box, halfspace, and no constraints

- [Moré & Wild, 2009], [Moré, Garbow, Hillstrom, 1981]



$$\tau = 10^{-1}$$



$$\tau = 10^{-3}$$

Constrained DFO: Summary

- ▶ Can ensure good geometry
 - ⇒ fully linear model
 - ⇒ error bound on approx. criticality measure
 - ⇒ convergence

- ▶ Worst-case complexity same as in unconstrained case

Constrained DFO: Future Work

- ▶ Convergence and WCC theory for quadratic models
- ▶ Fully quadratic model, etc.

References

- [1] Amir Beck. *First-Order Methods in Optimization*. Oct. 2017.
- [2] Coralia Cartis and Lindon Roberts. “A derivative-free Gauss-Newton method”. In: *Mathematical Programming Computation* 11.4 (2019), pp. 631–674.
- [3] Andrew R. Conn, Nicholas I. M. Gould, and Philippe L. Toint. *Trust-Region Methods*. Vol. 1. MPS-SIAM Series on Optimization. Philadelphia: MPS/SIAM, 2000.
- [4] Andrew R. Conn, Katya Scheinberg, and Luís N. Vicente. *Introduction to Derivative-Free Optimization*. Vol. 8. MPS-SIAM Series on Optimization. Philadelphia: MPS/SIAM, 2009.
- [5] Jorge J. Moré, Burton S. Garbow, and Kenneth E. Hillstom. “Testing Unconstrained Optimization Software”. In: *ACM Transactions on Mathematical Software* 7.1 (Mar. 1981), pp. 17–41.
- [6] Jorge J. Moré and Stefan M. Wild. “Benchmarking Derivative-Free Optimization Algorithms”. In: *SIAM Journal on Optimization* 20.1 (Jan. 2009), pp. 172–191.