

Derivative-Free Optimization with Convex Constraints

Joint work with Lindon Roberts (ANU → University of Sydney)

Matthew Hough, University of Waterloo (mough@uwaterloo.ca)

ICCOPT 2022, July 25-28, 2022

1. **Introduction to DFO trust-region methods**
2. Extending to convex constraints
3. Application to least-squares problems

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: $\|\nabla f(\mathbf{x}^*)\| \leq \epsilon$)

- 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:

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

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

Model-Based DFO: Algorithm

1. Build local interpolation model $m_k(\mathbf{s})$
2. Minimize the model within the trust-region Δ_k to get the step

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

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

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

4. **Update interpolation set:** add $\mathbf{x}_k + \mathbf{s}_k$ to the interpolation set
5. **If needed, ensure new interpolation set is 'good'**

You may be wondering...

1. What does it mean for our interpolation model to be a 'good approximation'?
2. What convergence/complexity guarantees do we have?

You may be wondering...

1. What does it mean for our interpolation model to be a 'good approximation'?
2. What convergence/complexity guarantees do we have?

An interpolation model $f(\mathbf{x}_k + \mathbf{s}) \approx m_k(\mathbf{s})$ is **fully linear** if

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

for all $\|\mathbf{s}\|_2 \leq \Delta_k$ (c.f. linear Taylor series).

Model-Based DFO: Theory

You may be wondering...

1. What does it mean for our interpolation model to be a 'good approximation'?
2. What convergence/complexity guarantees do we have?

An interpolation set is Λ -poised if

$$\max_t \max_{\|\mathbf{s}\|_2 \leq \Delta_k} |\ell_t(\mathbf{x}_k + \mathbf{s})| \leq \Lambda,$$

where $\ell_t(\mathbf{y}_s) = \delta_{s,t}$.

Theorem

If the interpolation set is Λ -poised and is contained in $B(\mathbf{x}_k, \Delta_k)$, then the corresponding interpolation model is fully linear with constants κ_{ef}, κ_{eg} in $\mathcal{O}(\Lambda)$. (+ dependencies on n, f)

Model-Based DFO: Theory

You may be wondering...

1. What does it mean for our interpolation model to be a 'good approximation'?
2. What convergence/complexity guarantees do we have?

Convergence and worst-case complexity for nonconvex functions match what we have for the derivative-based case.

Theorem (convergence)

Suppose f has Lipschitz continuous gradient and is bounded below. Then $\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_k)\|_2 = 0$.

Model-Based DFO: Theory

You may be wondering...

1. What does it mean for our interpolation model to be a 'good approximation'?
2. [What convergence/complexity guarantees do we have?](#)

Convergence and worst-case complexity for nonconvex functions match what we have for the derivative-based case.

Theorem (convergence)

Suppose f has Lipschitz continuous gradient and is bounded below. Then $\lim_{k \rightarrow \infty} \|\nabla f(\mathbf{x}_k)\|_2 = 0$.

Theorem (complexity)

Under the same assumptions as above, we achieve $\|\nabla f(\mathbf{x}_k)\|_2 \leq \epsilon$ for the first time after at most $\mathcal{O}(\epsilon^{-2})$ iterations. (+ dependencies on n, f)

1. Introduction to DFO trust-region methods
2. **Extending to convex constraints**
3. Application to least-squares problems

The 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 $\nabla f(\mathbf{x})$
- $\mathcal{C} \subseteq \mathbb{R}^n$ has nonempty interior, closed, and convex
 - Strictly feasible algorithm: never evaluate f at points outside \mathcal{C} ;
 - Access to \mathcal{C} is only through a (cheap) projection operator

Examples: \mathbb{R}^n , bound constraints, half-plane, Euclidean ball, ...

- Unrelaxable constraints
 - Only done for simple cases, no convergence theory
 - Bounds [Powell, 2009; Wild, 2009; Gratton et al., 2011]
 - Linear inequalities [Gumma, Hashim & Ali, 2014; Powell, 2015]
- Convex constraints with projection [Conejo et al., 2013]
 - Convergence, no complexity
 - Assume models are always fully linear
- Derivative-based complexity analysis [Cartis, Gould & Toint, 2012]

Existing work

- Unrelaxable constraints
 - Only done for simple cases, no convergence theory
 - Bounds [Powell, 2009; Wild, 2009; Gratton et al., 2011]
 - Linear inequalities [Gumma, Hashim & Ali, 2014; Powell, 2015]
- Convex constraints with projection [Conejo et al., 2013]
 - Convergence, no complexity
 - Assume models are always fully linear
- Derivative-based complexity analysis [Cartis, Gould & Toint, 2012]

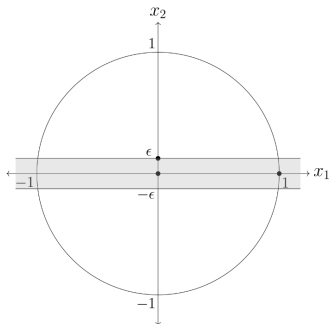
What is the problem? (Larson, Menickelly & Wild, 2019)

Model-based methods are more challenging to design in the presence of unrelaxable constraints because enforcing guarantees of model quality such as full linearity can be difficult. For fixed κ_{ef}, κ_{eg} it may be impossible to obtain a fully linear model using only feasible points.

What can we do?

Convex Constraints: The Problem

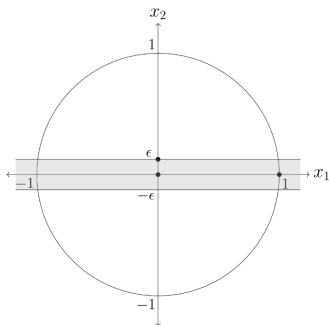
- $\mathcal{C} = \{(x_1, x_2) : |x_2| \leq \epsilon\} \subseteq \mathbb{R}^2$
- $Y = \{(0, 0), (1, 0), (0, \epsilon)\} \subseteq B(\mathbf{0}, 1)$



In $B(\mathbf{0}, 1)$, points are Λ -poised with $\Lambda = \mathcal{O}(\epsilon^{-1}) \implies$ large interpolation error. We cannot improve this using only feasible points.

Convex Constraints: The Problem

- $\mathcal{C} = \{(x_1, x_2) : |x_2| \leq \epsilon\} \subseteq \mathbb{R}^2$
- $Y = \{(0, 0), (1, 0), (0, \epsilon)\} \subseteq B(\mathbf{0}, 1)$



In $B(\mathbf{0}, 1)$, points are Λ -poised with $\Lambda = \mathcal{O}(\epsilon^{-1}) \implies$ large interpolation error. We cannot improve this using only feasible points.

If we only consider $|\ell_t(\mathbf{x}_k + \mathbf{s})|$ inside the feasible region $\implies \Lambda = \mathcal{O}(1)$.

Recall the old definition of a Λ -poised set:

$$\max_t \max_{\|\mathbf{s}\|_2 \leq \Delta_k} |\ell_t(\mathbf{x}_k + \mathbf{s})| \leq \Lambda,$$

Convex Constraints: Geometry

Recall the old definition of a Λ -poised set:

$$\max_t \max_{\|\mathbf{s}\|_2 \leq \Delta_k} |\ell_t(\mathbf{x}_k + \mathbf{s})| \leq \Lambda,$$

New definition:

$$\max_t \max_{\substack{\mathbf{x}_k + \mathbf{s} \in \mathcal{C} \\ \|\mathbf{s}\|_2 \leq \Delta_k}} |\ell_t(\mathbf{x}_k + \mathbf{s})| \leq \Lambda,$$

- We only care about the magnitude of Lagrange polynomials inside the feasible region.
- Gives smaller values of Λ .

Convex Constraints: Geometry

Recall the old definition of a fully linear model:

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

This is unnecessarily strong for our case!

Convex Constraints: Geometry

Recall the old definition of a fully linear model:

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

This is unnecessarily strong for our case!

New definition:

$$\begin{aligned} \max_{\substack{\mathbf{x}_k + \mathbf{s} \in \mathcal{C} \\ \|\mathbf{s}\|_2 \leq \Delta_k}} |f(\mathbf{x}_k + \mathbf{s}) - m_k(\mathbf{s})| &\leq \kappa_{ef} \Delta_k^2, \\ \max_{\substack{\mathbf{x}_k + \mathbf{s} \in \mathcal{C} \\ \|\mathbf{s}\|_2 \leq 1}} \|(\nabla f(\mathbf{x}_k) - \nabla m_k(\mathbf{0}))^T \mathbf{s}\|_2 &\leq \kappa_{eg} \Delta_k \end{aligned}$$

Theorem (Hough & Roberts, 2021)

If the set of interpolation points is contained in $B(\mathbf{x}_k, \Delta_k) \cap \mathcal{C}$ and is Λ -poised, then the corresponding **linear** interpolation model is fully linear with $\kappa_{ef}, \kappa_{eg} = \mathcal{O}(\Lambda)$.

Convex Constraints: The Algorithm

1. Build local interpolation model $m_k(\mathbf{s})$
2. Minimize the model within the trust-region Δ_k to get the step

$$\mathbf{s}_k = \arg \min_{\mathbf{s} \in \mathbb{R}^n} m_k(\mathbf{s}) \quad \text{s.t.} \quad \|\mathbf{s}\|_2 \leq \Delta_k \quad \text{and} \quad \mathbf{x}_k + \mathbf{s} \in \mathcal{C}$$

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

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

4. Update interpolation set: add $\mathbf{x}_k + \mathbf{s}_k$ to the interpolation set
5. If needed, ensure new interpolation set is Λ -poised

Convex Constraints: Convergence & Complexity

A new question arises...

How can we measure stationarity? i.e. how do we measure progress?

Convex Constraints: Convergence & Complexity

A new question arises...

How can we measure stationarity? i.e. how do we measure progress?

$$\pi^f(\mathbf{x}) := \left| \min_{\substack{\mathbf{x}+\mathbf{s} \in \mathcal{C} \\ \|\mathbf{s}_2\| \leq 1}} \nabla f(\mathbf{x})^T \mathbf{s} \right|$$

Properties:

[Conn, Gould & Toint, 2000]

- $\pi^f(\mathbf{x}) \geq 0$ for all \mathbf{x}
- $\pi^f(\mathbf{x}^*) = 0$ if and only if \mathbf{x}^* is a KKT point
- If $\mathcal{C} = \mathbb{R}^n$, then $\pi^f(\mathbf{x}) = \|\nabla f(\mathbf{x})\|_2$
- $\pi^f(\mathbf{x})$ is Lipschitz continuous in \mathbf{x} , assuming ∇f is Lipschitz continuous [Cartis, Gould & Toint, 2012]

Convex Constraints: Convergence & Complexity

A new question arises...

How can we measure stationarity? i.e. how do we measure progress?

$$\pi^f(\mathbf{x}) := \left| \min_{\substack{\mathbf{x}+\mathbf{s} \in \mathcal{C} \\ \|\mathbf{s}_2\| \leq 1}} \nabla f(\mathbf{x})^T \mathbf{s} \right|$$

Properties:

[Conn, Gould & Toint, 2000]

- $\pi^f(\mathbf{x}) \geq 0$ for all \mathbf{x}
- $\pi^f(\mathbf{x}^*) = 0$ if and only if \mathbf{x}^* is a KKT point
- If $\mathcal{C} = \mathbb{R}^n$, then $\pi^f(\mathbf{x}) = \|\nabla f(\mathbf{x})\|_2$
- $\pi^f(\mathbf{x})$ is Lipschitz continuous in \mathbf{x} , assuming ∇f is Lipschitz continuous [Cartis, Gould & Toint, 2012]
- If m_k is fully linear, then $|\pi^f(\mathbf{x}_k) - \pi^{m_k}(\mathbf{x}_k)| \leq \kappa_{eg} \Delta_k$ [Hough & Roberts, 2021]

Convex Constraints: Convergence & Complexity

Convergence and worst-case complexity for nonconvex functions match what we have for the unconstrained case:

Theorem (convergence) (Hough & Roberts, 2021)

If f has Lipschitz continuous gradient and is bounded below, then we have $\lim_{k \rightarrow \infty} \pi^f(\mathbf{x}_k) = 0$.

Convex Constraints: Convergence & Complexity

Convergence and worst-case complexity for nonconvex functions match what we have for the unconstrained case:

Theorem (convergence) (Hough & Roberts, 2021)

If f has Lipschitz continuous gradient and is bounded below, then we have $\lim_{k \rightarrow \infty} \pi^f(\mathbf{x}_k) = 0$.

Theorem (complexity) (Hough & Roberts, 2021)

Under the same assumptions as above, we achieve $\pi^f(\mathbf{x}_k) \leq \epsilon$ for the first time after at most $\mathcal{O}(\epsilon^{-2})$ iterations.

Convex Constraints: Convergence & Complexity

Convergence and worst-case complexity for nonconvex functions match what we have for the unconstrained case:

Theorem (convergence) (Hough & Roberts, 2021)

If f has Lipschitz continuous gradient and is bounded below, then we have $\lim_{k \rightarrow \infty} \pi^f(\mathbf{x}_k) = 0$.

Theorem (complexity) (Hough & Roberts, 2021)

Under the same assumptions as above, we achieve $\pi^f(\mathbf{x}_k) \leq \epsilon$ for the first time after at most $\mathcal{O}(\epsilon^{-2})$ iterations.

Requires two important algorithms:

- Check a model is fully linear
- Change the interpolation set to make the model fully linear if it is not

1. Introduction to DFO trust-region methods
2. Extending to convex constraints
3. **Application to least-squares problems**

The Least-Squares Case

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

- Typically (Gauss-Newton) linearize \mathbf{r} at \mathbf{x}_k using the Jacobian:

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

- But in DFO, Jacobian is not available:

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

- Find \mathbf{J}_k by interpolation [Cartis & Roberts, 2019]

The Least-Squares Case

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

- Typically (Gauss-Newton) linearize \mathbf{r} at \mathbf{x}_k using the Jacobian:

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

- But in DFO, Jacobian is not available:

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

- Find \mathbf{J}_k by interpolation [Cartis & Roberts, 2019]

Either way, end up with a local quadratic model

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

Least-Squares Implementation

Update the state-of-the-art solver **DFO-LS**: [Cartis et al., 2019]

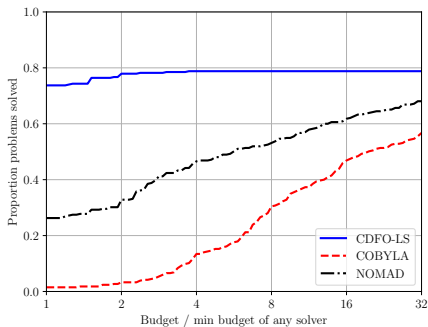
- Use FISTA to solve the constrained trust-region subproblem
- Requires Dykstra's algorithm to project onto $B(\mathbf{x}_k, \Delta_k) \cap \mathcal{C}$
- Github: [numericalalgorithmsgroup/dfols](https://github.com/numericalalgorithmsgroup/dfols)

Tested on a collection of 58 low-dimensional least-squares problems with box/ball/halfspace/second-order cone constraints.

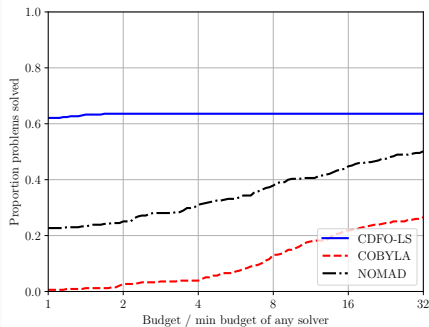
Limited solvers to compare to (none exploit the least-squares structure):

- NOMAD: direct search DFO, models constraints using a barrier method (i.e. $f(\mathbf{x}) = +\infty$ if $\mathbf{x} \notin \mathcal{C}$) [Le Digabel, 2011]
- COBYLA: model-based DFO with inequality constraints [Powell, 1994]

Numerical Results



Low accuracy, $\tau = 10^{-1}$



High accuracy, $\tau = 10^{-5}$

Measuring the proportion of problems solved vs. the number of objective evaluations (higher is better).

Summary & Future Work

Summary

- General model-based DFO method for convex-constrained problems
- Match/generalize existing convergence and complexity results
- Develop new theory of Λ -poisedness and full linearity¹
- New software for least-squares problems

Future Work

- Second-order theory
- Generalize interpolation theory to quadratic interpolation

[[arXiv:2111.05443](https://arxiv.org/abs/2111.05443), [Github: numericalalgorithmsgroup/dfols](https://github.com/numericalalgorithmsgroup/dfols)]

¹Only for linear/composite models at present