

Possible Future Research

Trever Hallock

April 13, 2017

Contents

1	Introduction	1
2	Derivative-free Mentality	2
2.0.1	Derivatives are convenient	2
2.0.2	Problem statement	2
2.1	Approaches	2
2.1.1	Finite difference methods	2
2.1.2	Direct search methods	3
2.1.3	Model based methods	3
2.2	The Derivative Trust Region Method	3
2.3	More details	3
3	Literature Review	3
4	Possible future work	4
4.1	Future Directions	4
4.2	NLP Methods	4
4.3	Goals	5
4.4	Parallelization	5
4.5	Miscellaneous	5
5	My Algorithm	5
5.1	Problem	5
5.2	First approach	5
5.3	Algorithm in other paper	5
5.3.1	Criticality Measure	6
5.3.2	Step decomposition	6
5.3.3	The filter	6
5.3.4	f -type versus θ -type	7
5.3.5	Restoration Step	7
5.4	When we use derivative-free methods	7
5.4.1	The change in the criticality measure	7
5.5	Algorithm Description	7
5.5.1	Simplified version draft	7
5.5.2	Psuedo code	7
5.6	Generalizations	9
6	Pictures of convergence	9
7	Comparison to other libraries (If I get enough time)	9

1 Introduction

This paper will discuss research topics for derivative free optimization. It begins with an introduction to the context and goals of derivative free optimization (DFO) supplemented by some of the recent advancement made in the field. It then details several future research directions and one in particular (filter methods) that has been studied. The focus is on local search algorithms for constrained derivative free optimization.

2 Derivative-free Mentality

Derivative free optimization is motivated by programs in which derivative information is unknown, deceptive or otherwise impractical to compute. For example, this can arise when the objective is the result of a simulation that does not admit automatic differentiation due to copyrights or is too expensive for repeated calls within a finite difference framework. The general strategy is to converge to a first or second order critical point while evaluating the function as few times as possible.

Another branch of DFO is concerned with noisy function evaluation. Noisy functions can be categorized as either deterministic or random. Deterministic means that the genuine objective is not always evaluated accurately, but the measurement error (?) will not change across multiple function calls at the same point. Random means that each point in the domain is associated with a distribution of possible values the objective may return.

2.0.1 Derivatives are convenient

The lack of derivative information means that DFO methods are at a disadvantage when compared to their counterparts in nonlinear optimization. First and second derivative information is explicit in algorithms with quadratic convergence such as Newton’s method. They are also present in conditions for convergence results such as Wolf’s, Armijo or Goldstien for line search methods. Additionally, stopping criteria usually involve a criticality test involving derivatives.

2.0.2 Problem statement

In more detail, DFO methods consider nonlinear, constrained optimization problems of the form

$$\begin{aligned} & \underset{x}{\text{minimize}} && f(x, S(x)) \\ & \text{subject to} && g_i(x, S(x)) \leq 0, \quad i = 1, \dots, m_{\mathcal{I}} \\ & && h_i(x, S(x)) = 0, \quad i = 1, \dots, m_{\mathcal{E}} \end{aligned}$$

which give rise to different classes of derivative-free optimization based on properties of $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $g_i : \mathbb{R} \rightarrow \mathbb{R}$, and $h_i : \mathbb{R} \rightarrow \mathbb{R}$. As discussed, these functions may be noisy, either deterministically or randomly, although this paper does not deal with this case. Derivatives of g_i and h_i may be known, in which case the objective is the only derivative free function. One common such case is to include bound constraints of the form $b_L \leq x \leq b_U$ for some $b_L < b_U$, which gives rise to Box Constrained DFO (BCDFO). If the constraints can be evaluated at points outside the feasible region, the constraints are called relaxable constraints. We consider the case where no derivative information of g_i and h_i are known, for example if they are also output from a simulation used to evaluate the objective. This means that each call to the objective gives values of the constraints as well, and vice versa. However, although derivatives are not known, we do assume that f , g_i , and h_i are all continuously twice differentiable. Some problems additionally contain “hidden” constraints which are not explicit in the model but merely result in a notification that the objective could not be evaluated at the requested point. This may mean that it is not possible to tell how close to a “hidden” constraint the current iterate lies. Again, we only consider only local search algorithms for this problem that seek a first or second order stationary point.

Many DFO methods simply let $f(x, S(x)) = S(x)$.

A slightly more general form is presented by [2] and given here:

$$\begin{aligned} & \underset{x}{\text{minimize}} && f(x, y, z) \\ & \text{subject to} && g_i(x, y, z) \leq 0 \quad \forall i = 1, \dots, m_1 \\ & && h_i(x, y, z) = 0, \quad \forall i = 1 \dots, m_2 \\ & && y = d(z) \end{aligned}$$

Biegler assumes that the dimension of u is small compared to the dimension of x and z .

2.1 Approaches

In this section we discuss several of the approaches taken when derivative information is not known.

2.1.1 Finite difference methods

Finite difference methods can be used to approximate the derivative of a function f . A common approximation is given by $\nabla f(x) \approx \left(\frac{f(x+he_i) - f(x-he_i)}{2h} \right)_i$ for some small h . This may work well but can have issues with unlucky iterates. The number of function evaluations tends to grow large with the dimension and number of iterations the algorithm performs as information is only gathered near the current iterate when h is small (which is required for accurate derivatives). Because of the large number of function evaluations required for finite difference schemes, these can be unusable.

CITE

2.1.2 Direct search methods

Another approach is to use direct search methods that do not explicitly estimate the derivative but evaluate the objective on a pattern or other structure to find a descent direction. Examples of this include Coordinate descent and other pattern based search methods. One of the most popular direct search method is Nelder Mead (it is implemented in `fminsearch` in matlab) although it is proven to not converge in (synonym: bad) cases. These methods can be robust but ignore information because they do not use derivatives.

(0th derivative)

2.1.3 Model based methods

In this paper, we are concerned with model based methods that minimize a model used to approximation of the unknown objective (by regression or kriging or ...). The algorithm can then use derivative information from the model to find a descent direction. These methods require an adequate geometry of sampled points to ensure convergence.

2.2 The Derivative Trust Region Method

- Interpolate or regress model functions onto a sample
- Minimize the model function over a trust region
- Adjust trust region, possibly by testing the model's accuracy

2.3 More details

1. Define $m_k(x) = f(x^{(k)}) + \nabla f(x^{(k)})^T(x - x^{(k)}) + \frac{1}{2}(x - x^{(k)})^T \nabla^2 f(x^{(k)})(x - x^{(k)})$
 - $\nabla f(x^{(k)})$ and $\nabla^2 f(x^{(k)})$ must be approximated
 - There are geometric properties of the sample set that must be satisfied
2. if $\nabla m_k(x) < t$ stop
3. Solve the Trust region subproblem: $s^{(k)} = \arg \min_{s \in B(x^{(k)}, \Delta_k)} m_k(x^{(k)} + s)$
4. Test for improvement
 - $\rho_k = \frac{f(x^{(k)}) - f(x^{(k)} + s^{(k)})}{m_k(x^{(k)}) - m_k(x^{(k)} + s^{(k)})}$
 - If ρ is small, $x^{(k+1)} = x^{(k)}$ (reject) and decrease radius
 - If ρ is intermediate, $x^{(k+1)} = x^{(k)} + s^{(k)}$ (accept) and decrease radius
 - If ρ is large, $x^{(k+1)} = x^{(k)} + s^{(k)}$ (accept) and increase radius

3 Literature Review

The original filter method was proposed by Gould in [1]. The motivation for the filter method was that the algorithm does not need to tune any parameters as in penalty or merit methods.

A recent paper from September 2016 [2] implements a derivative-free trust region filter method for solving the general program

$$\begin{aligned}
 & \underset{x}{\text{minimize}} && f(x, y, z) \\
 & \text{subject to} && g_i(x, y, z) \leq 0 \quad \forall i = 1, \dots, m_1 \\
 & && h_i(x, y, z) = 0, \quad \forall i = 1, \dots, m_2 \\
 & && y = d(z)
 \end{aligned}$$

where d is a black box function while f , g and h are glass box functions. This is more general than the algorithm we consider as it allows for the objective to depend on other glass box functions of the input and multiple outputs of the black box function. The authors compare their algorithm to finite difference methods as well as kriging on three different applications from Chemistry. Within the algorithm, the current iterate is always feasible with respect all inequalities except for the constraint $y = d(z)$.

In the 2006 paper [3] Fletcher reviews how filter methods have developed. The only references to derivative-free versions of the filter method are those applied to pattern based (direct) methods. However, the authors outline trust region filter methods along with several other variants.

In Brekelman’s paper [4], a trust region filter method is developed to minimize function evaluations by constructing linear model functions. This paper also uses experimental designs to choose following iterates.

Within [5] and [6] Biegler uses a filter method to ensure global convergence within a line search framework. We experimented with this before deciding combining line search with the derivative-free trust region approach was not a natural approach. (There have been attempts to employ both of these frameworks in [7].)

Within [8] derivative-free methods are developed in detail. This contains a good explanation of ensuring geometry of the current set with poisedness for unconstrained problems and also covers other derivative-free methods including direct-search and line search.

Within [9] Toint generalizes the filter method with the notion of a trust funnel. This is for glass box functions as it does not include any derivative-free methods.

Within [10] a sequential quadratic programming method is applied to the filter method?

Colson has also applied filter techniques to derivative-free optimization in his 2004 Ph.D. thesis [11]. This focuses on bilevel programming.

I based my algorithm of the trust region filter method described in:

GLOBAL CONVERGENCE OF A TRUST-REGION SQP-FILTER ALGORITHM FOR GENERAL NONLINEAR PROGRAMMING ROGER FLETCHER, NICHOLAS I. M. GOULD, SVEN LEYFFER, PHILIPPE L. TOINT, AND ANDREAS WACHTER I need to find a bibtex reference of this paper to put in the bibliography.

the two additional review papers.

4 Possible future work

4.1 Future Directions

- Convert classical NLP Algorithms to DFO
- Different Goals
- Parallelization
- Explore different model functions
- Providing structure to the optimization program

4.2 NLP Methods

We can convert classical algorithms to DFO Possible future work includes converting nonlinear algorithms to a derivative free context. Several algorithms have already been converted.

- Line search (CITE)
- Filter method (introduced in 2004 CITE for pattern searches)
- Active Set (CITE)
- Augmented Lagrangian (CITE)
- Penalty
 - Abramson & Audet, 2006
 - Abramson et al. 2009c
 - audet et al. 2008b
 - Audet & Dennis, 2006
 - sequential penalty merit functions Liuzzi et al., 2010
 - smoothed exact $l - \infty$ penalty function Liuzzi & Lucidi, 2009
 - exact penalty merit function Fasano, Liuzzi, Lucidi, & Rinaldi, 2014; Gratton & Vicente, 2014
- Progressive Barrier
 - Audet & Dennis, 2009
- Interior Point (CITE)

The approach is to replace derivatives of f , g_i , h_i with derivatives of a model function. I have considered line search methods as well as Filter methods which are discussed below.

4.3 Goals

We could consider an objective function with a evaluation runtime that varies with x .

4.4 Parallelization

1. parallelization of the function and/or the derivative evaluations in the algorithm
2. parallelization of linear algebra kernels
3. modifications of the basic algorithms which increase the degree of intrinsic parallelism, for instance, by performing multiple function and/or derivative evaluations

The third item is particularly interesting.

4.5 Miscellaneous

- Model Functions Radial basis functions
- Problem Structure Specify the structure of f , g_i , or h_i
- Machine Learning Applying DFO techniques to algorithms that currently use random sampling to tune parameters

5 My Algorithm

5.1 Problem

We consider problems of the form

$$\begin{aligned} \min_x f(x) \\ g(x) &\leq 0 \\ h(x) &= 0 \end{aligned}$$

where $f : R^n \rightarrow R$, $g : R^n \rightarrow R^{m_1}$ and $h : R^n \rightarrow R^{m_2}$. We assume that all functions are derivative-free: all functions f, g, h are evaluated by a single call to a black box function $d(x) = (f(x), g(x)^T, h(x)^T)^T$.

5.2 First approach

We began with a line search filter method, however we found that this had several drawbacks:

- Trust regions arise naturally within derivative-free algorithms
- Line search algorithms exploit how much easier finding a descent direction is than solving the trust region subproblem, but this saved computation is not as useful in DFO
- As the algorithm backtracks on the step length, the trust region must be reduced, as the models are accurate over a region rather than at a single point

5.3 Algorithm in other paper

We work within a trust region, sequential quadratic programming framework that uses a filter method introduced by Fletcher.

We first compute an interpolation set poised for regressing a set of model functions, which we choose to be quadratic functions. Although we have function values for enough points to construct model functions of the same order as used for the objective, we model the feasible with only the linear constraints.

The core of the algorithm revolves around

5.3.1 Criticality Measure

In order to construct stopping criteria, we introduce a criticality measure χ which goes to zero as the iterates approach a first order critical point.

This is defined as

$$\chi = \begin{array}{ll} |\min_t \langle g_k + H_k n_k, t \rangle| & \\ A_{eq} t & = 0 \\ c_{ineq} + A_{ineq} t & \leq 0 \\ \|t\| & \leq 1 \end{array}$$

5.3.2 Step decomposition

At iteration k , we can decompose the step s_k into a normal step n_k intended to decrease constraint violation and a tangential step t_k intended to reduce the objective. The step n_k projects the current iterate onto the feasible region. Currently, we project x_k onto only the linear model of the feasible region. We require that the the computation of the normal step, which solves:

$$\begin{array}{ll} n_k = & \arg \min_n \|n\|^2 \\ s.t. & c_{eq} + A_{eq} n = 0 \\ & c_{ineq} + A_{ineq} n \leq 0 \\ & \|n\|^2 \leq \Delta_k^2 \end{array}$$

In addition to this program having a feasible point, we need to know that there is enough space for us to provide sufficient decrease within the tangential step. This means that we require the stronger condition that

$$\|n\| \leq \kappa_\Delta \Delta_k \min\{1, \kappa_\mu \Delta_k^\mu\}$$

If this condition is satisfied, then we say that the program is *compatible*. We are then able to compute a tangential step t_k :

$$\begin{array}{ll} t_k = & \arg \min_t (g_n + H_k n_k)^T t + \frac{1}{2} t^T H_k t \\ s.t. & c_{eq} + A_{eq} t = 0 \\ & c_{ineq} + A_{ineq} t \leq 0 \\ & \|n_k + t_k\|^2 \leq \Delta_k^2 \end{array}$$

- quadratic information contained in H_k
- H_k is the hessian of the lagrangian, as computed by using KKT the matrix
- This program is a shifted version of another

5.3.3 The filter

The filter is a method used to ensure convergence to a feasible point. It works by ensuring that all new iterates are nondominated with respect to all previous iterates when the problem is viewed as a multi-criteria optimization problem $\min(\theta, f)$. However, simply ensuring that new points are nondominated does not provide sufficient progress, we must ensure that the objective decreases by a greater amount when the current iterate is far from the feasible region:

$$\theta(x_k) \leq (1 - \gamma_\theta) \theta_i$$

or

$$f(x_k) \leq f_i - \gamma_\theta \theta_i$$

for all (θ_i, f_i) that are currently in the filter. When this condition is satisfied, we say that the new iterate x_k is admissible to the filter.

5.3.4 f -type versus θ -type

When steps decrease f significantly more than θ the steps are considered f type. θ -type steps are steps that significantly reduce the constraint violation. I thought that the paper I read about the convergence rate said they introduced a new criteria to this that made it more efficient.

5.3.5 Restoration Step

The goal of the feasibility restoration step is to find a new iterate and trust region radius that allows the current iterate to be compatible.

While performing the restoration step, we place constraints in the objective by minimizing the squared norm of θ . We take one step to minimize the quadratic model unconstrained optimization problem, and update the trust region based on the new function value.

It is possible that the restoration step is unsuccessful if the iterates approach an infeasible local minimum of the constraints. In this case, the algorithm will fail to find a feasible local minimum, and will need to be restarted.

5.4 When we use derivative-free methods

5.4.1 The change in the criticality measure

One issue with applying the original algorithm within a DFO context was that the trust region radius is not required to go to zero. However, within DFO we must also require that the trust region goes to zero as we approach a stationary point. One way of ensuring this is to decrease the trust region radius when the current step lies well within the trust region radius. However, a better approach may be to introduce a tolerance on the criticality measure and decrease the trust region whenever the criticality falls below the threshold.

5.5 Algorithm Description

5.5.1 Simplified version draft

- compute model functions and hessian of the lagrangian
- compute criticality measure
- if feasible and critical, then decrease trust region radius or return
- compute normal step
- check compatibility, restoring feasibility if necessary and return to step 1
- compute tangential step
- check for sufficient reduction in the model function, adding the current iterate to the filter and returning to step 1 if necessary
- evaluate the function and constraints at the trial point
- check compatibility to the filter, decreasing the trust region radius and returning to step 1 if necessary
- compute ρ , and accept the trial point or reject and decrease the radius

5.5.2 Psuedo code

Algorithm 1 Filter Trust Region Search

```

1: procedure TRUST REGION FILTER
2:   initialize
3:    $k = 0$ 
4:   choose an  $x_0$ 
5:   while  $k < \text{maxit}$  do
6: main loop:
7:   ensure poisedness, possibly adding points to the model
8:   Compute  $m_k, g_k = \nabla m_k(x_k), c_k, A_k, f_k = f(x_k), \mathcal{A}, \theta_k$ 
9:   Solve:
10:      
$$\begin{aligned} \nabla^2 m_k(x_k) d + A_k^T \lambda &= g_k \\ A_k d &= c_k \end{aligned}$$

11:    $H_k \leftarrow \nabla^2 m_k(x_k) + \sum_i \lambda_i \nabla^2 c_{ik}$ 
12:    $\chi_k \leftarrow |\min_t \{ \langle g_k + H_k n_k, t \rangle | A_{eq} t = 0 \wedge c_{ineq} + A_{ineq} t \leq 0 \wedge \|t\| \leq 1 \}|$ 
13:   if constraint violation  $= 0 \wedge \chi = 0$  then
14:     if  $\text{tol} < \Delta_k$  then
15:       reduce  $\Delta$ :  $\Delta_{k+1} \leftarrow \text{some} \in [\gamma_0 \Delta_k, \gamma_1 \Delta_k]$ 
16:        $k \leftarrow k + 1$ 
17:       go to main loop
18:     success
19:      $n_k \leftarrow \arg \min_n \{ \|n\|^2 | c_{eq} + A_{eq} n = 0 \wedge c_{ineq} + A_{ineq} n \leq 0 \wedge \|n\|^2 \leq \Delta_k \}^2$ 
20:     if Feasible region  $\neq \emptyset \wedge \|n\| \leq \kappa_\Delta \Delta_k \min\{1, \kappa_\mu \Delta_k^\mu\}$  then
21:        $t_k \leftarrow \arg \min_t \{ (g_n + H_k n_k)^T t + \frac{1}{2} t^T H_k t | c_{eq} + A_{eq} t = 0 \wedge c_{ineq} + A_{ineq} t \leq 0 \wedge \|s\| \leq \Delta_k \}$ 
22:        $s_k \leftarrow t_k + n_k$ 
23:       if  $m_k(x_k) - m_k(x_k + s_k) \geq \kappa_\theta \theta_k^\psi$  then
24:         add  $x_k$  to filter
25:         reduce  $\Delta$ :  $\Delta_{k+1} \leftarrow \text{some} \in [\gamma_0 \Delta_k, \gamma_1 \Delta_k]$ 
26:         go to main loop
27:       // Here we evaluate new  $c$  and  $f$  at  $x_k + s_k$ 
28:       if  $x_k + s_k$  is acceptable:  $\theta(x_k + s_k) \leq (1 - \gamma_\theta) \theta' \vee f(x_k + s_k) \leq f' - \gamma_\theta \theta' \forall (f', \theta') \in \text{Filter}$  then
29:          $\rho = \frac{f(x_k) - f(x_k + s_k)}{m_k(x_k) - m_k(x_k + s_k)}$ 
30:         if  $\rho < \eta_1$  then
31:           reduce  $\Delta$ :  $\Delta_{k+1} \leftarrow \text{some} \in [\gamma_0 \Delta_k, \gamma_1 \Delta_k]$ 
32:            $k \leftarrow k + 1$ 
33:           go to main loop
34:         else if  $\rho > \eta_2$  then
35:           if  $\|s\| < \frac{\Delta_k}{2}$  then
36:             reduce  $\Delta$ :  $\Delta_{k+1} \leftarrow \text{some} \in [\gamma_0 \Delta_k, \gamma_1 \Delta_k]$ 
37:           else
38:             increase  $\Delta$ :  $\Delta_{k+1} \leftarrow \text{some} \in [\Delta_k, \gamma_2 \Delta_k]$ 
39:            $x_{k+1} \leftarrow x_k + s_k$ 
40:         else
41:           reduce  $\Delta$ :  $\Delta_{k+1} \leftarrow \text{some} \in [\gamma_0 \Delta_k, \gamma_1 \Delta_k]$ 
42:            $k \leftarrow k + 1$ 
43:           go to main loop
44:         else
45:           add  $x_k$  to filter
46:           compute new  $r$  (restoration step) and  $\Delta$ 
47:           if impossible to restore then fail
48:            $x_{k+1} \leftarrow x_k + r$ 
49:            $k \leftarrow k + 1$ 

```

5.6 Generalizations

As this method was also discovered by () we have also considered possible extensions.

- The authors only considered linear models, we can use quadratic models
- We can generalize (relax) the possible steps by only requiring only that new iterates are admissible with respect to a multi objective filter of the form $(f(x), \|g_1(x)\|^2, \dots, \|g_{m_1}(x)\|^2, \|h_1(x)\|, \dots, \|h_{m_2}\|)$.
- We could relax the condition that all iterates have to remain feasible with respect to all constraints except $y = d(z)$.

6 Pictures of convergence

7 Comparison to other libraries (If I get enough time)

References

- [1] N. I. M. Gould and P. L. Toint, “Nonlinear programming without a penalty function or a filter,” *Mathematical Programming*, vol. 122, no. 1, pp. 155–196, 2010.
- [2] J. P. Eason and L. T. Biegler, “A trust region filter method for glass box/black box optimization,” *AIChE Journal*, vol. 62, no. 9, pp. 3124–3136, 2016.
- [3] R. Fletcher, S. Leyffer, R. Fletcher, and S. Leyffer, “A brief history of filter methods,” tech. rep., 2006.
- [4] R. Brekelmans, L. Driessen, H. Hamers, and D. den Hertog, “Constrained optimization involving expensive function evaluations: A sequential approach,” *European Journal of Operational Research*, vol. 160, no. 1, pp. 121 – 138, 2005. Applications of Mathematical Programming Models.
- [5] A. Wächter and L. T. Biegler, “Line search filter methods for nonlinear programming: Motivation and global convergence,” *SIAM Journal on Optimization*, vol. 16, no. 1, pp. 1–31, 2005.
- [6] A. Wächter and L. T. Biegler, “Line search filter methods for nonlinear programming: Local convergence,” *SIAM Journal on Optimization*, vol. 16, no. 1, pp. 32–48, 2005.
- [7] J. Nocedal and Y. Yuan, “Combining trust region and line search techniques,”
- [8] A. R. Conn, K. Scheinberg, and L. N. Vicente, *Introduction to Derivative-Free Optimization*. Philadelphia, PA, USA: Society for Industrial and Applied Mathematics, 2009.
- [9] R. Sampaio and L. Toint, “A trust-funnel method for nonlinear optimization problems with general nonlinear constraints and its application to derivative-free optimization, most information in this cite is wrong,” *NAXYS Namur Center for Complex Systems*, vol. 61, no. 5000, pp. 1–31, 2015.
- [10] P. T. Boggs and J. W. Tolle, “Sequential quadratic programming for large-scale nonlinear optimization,” *Journal of Computational and Applied Mathematics*, vol. 124, no. 12, pp. 123 – 137, 2000. Numerical Analysis 2000. Vol. IV: Optimization and Nonlinear Equations.
- [11] B. Colson, “Trust-region algorithms for derivative-free optimization and nonlinear bilevel programming,” *Quarterly Journal of the Belgian, French and Italian Operations Research Societies*, vol. 2, no. 1, pp. 85–88, 2004.