

NEW YORK UNIVERSITY

MASTER'S THESIS

An L-BFGS-B Optimizer for Non-Smooth Functions

Author:

Wilmer Henao

Supervisor:

Michael L. Overton

*A thesis submitted in fulfilment of the requirements
for the degree of Master of Science in Scientific Computing*

in the

Courant Institute of Mathematical Sciences
Department of Mathematics

March 2014

Declaration of Authorship

I, Wilmer Henao, declare that this thesis titled, 'An L-BFGS-B Optimizer for Non-Smooth Functions' and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

NEW YORK UNIVERSITY

Abstract

Faculty Name

Department of Mathematics

Master of Science in Scientific Computing

An L-BFGS-B Optimizer for Non-Smooth Functions

by Wilmer Henao

The Thesis Abstract is written here ...

Acknowledgements

I would like to thank my advisor Michael Overton for all the hours of hard work and for all the great recommendations and changes that led to this thesis.

I also would like to thank Allan Kaku and Anders Skaaja for earlier contributions to other versions of the code.

Finally, I would like to thank High Performance Computing at NYU.

Contents

Declaration of Authorship	i
Abstract	ii
Acknowledgements	iii
Contents	iv
1 Introduction	1
2 L-BFGS-B	3
2.1 BFGS	3
2.2 L-BFGS	4
2.3 L-BFGS-B	4
2.3.1 Gradient Projection	4
2.3.2 Subspace Minimization	5
3 Modifications to the L-BFGS-B algorithm	7
3.1 The Wolfe conditions	7
3.2 The line search methodology	8
3.3 The termination condition	9
3.3.1 Description of the Solution	9
3.3.2 The Solution of the Quadratic Program	10
4 the solution test functions	12
A Appendix Title Here	14
Bibliography	15

LAH List Abbreviations **H**ere

For/Dedicated to/To my...

Chapter 1

Introduction

The problem addressed is to find a local minimizer of the nonsmooth minimization problem

$$\begin{aligned} \min_{x \in \mathbb{R}^n} \quad & f(x) \\ \text{s.t.} \quad & l_i \leq x_i \leq u_i, \\ & i = 1, \dots, n. \end{aligned} \tag{1.1}$$

where $f: \mathbb{R}^n \rightarrow \mathbb{R}$, is continuous but not differentiable everywhere and n is a very large but finite number.

The $L - BFGS - B$ algorithm [1] is a standard method for solving large instances of 1.1 when f is a smooth function. The original name of $BFGS$ stands for Broyden, Fletcher, Goldfarb and Shanno, the authors of the original "BFGS" quasi-Newton algorithm for unconstrained optimization discovered and published independently by them in 1970 [give the 4 references here]. This method requires storing and updating a matrix which approximates the inverse of the Hessian matrix $\nabla^2 f(x)$ and hence requires $\mathcal{O}(n^2)$ operations per iteration. The $L - BFGS$ variant [give reference] is based on BFGS but requires only $\mathcal{O}(mn)$ operations per iteration: thus, the L stands for Large. Finally, the last letter B in $L - BFGS - B$ stands for bounds, meaning the lower and upper bounds l_i and u_i in 1.1. The $L - BFGS - B$ algorithm is implemented in a well known *FORTRAN* software package by the same name [give reference].

In this thesis, there is a brief description of the $L - BFGS - B$ algorithm at a high level and then explain how the modified algorithm is more suitable for functions f which may not be differentiable at their local or global optimizers. We call the new algorithm L-BFGS-B-NS where NS stands for Non-Smooth. We implemented these changes in a

modified version of the Fortran code [ref] which is can be downloaded from the website for this thesis [give URL]. We report on some numerical experiments that strongly suggest that the new code should be useful for the nonsmooth bound-constrained optimization problem (1.1).

We are grateful to Jorge Nocedal and his coauthors for allowing us to modify the L-BFGS-B code and post the modified version.

Chapter 2

L-BFGS-B

This section is a description of the original $L - BFGS - B$ code at a very high level [2]. The original software is intended to work well with smooth functions. This thesis discusses how to modify the algorithm for Non-Smooth functions.

2.1 BFGS

BFGS is a standard tool for optimization of smooth functions.[3] It is a line search method and its goal is to find a search direction starting from its current position x . The search direction is of type $d = -B\nabla f$ ¹ where B is an approximation to the inverse Hessian. ² This k^{th} step approximation is calculated via the *BFGS* formula

$$B_{k+1} = \left(I - \frac{s_k y_k^T}{y_k^T s_k} \right) B_k \left(I - \frac{y_k s_k^T}{y_k^T s_k} \right) + \frac{s_k s_k^T}{y_k^T s_k} \quad (2.1)$$

where $y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$ and $s_k = x_{k+1} - x_k$ and where the first approximation of B_0 is assumed to be the identity matrix I in this thesis. *BFGS* exhibits superlinear convergence but it also requires $\mathcal{O}(n^2)$ operations per iteration. [3]

In the case of Non-Smooth functions. *BFGS* typically succeeds in finding a local minimizer. This however requires some modifications of the line search conditions. This line search conditions is known as the weak Wolfe line search and it will be explained later in this thesis.

¹Notice that when B is the identity, this is the same direction as steepest descent. Another common line search method of optimization

²When it is exactly the inverse Hessian the method is known as Newton's method. Newton's method has quadratic convergence but requires the explicit calculation of the Hessian at every single step.

2.2 L-BFGS

$L - BFGS$ stands for Limited-memory $BFGS$. This algorithm approximates $BFGS$ using only a limited amount of computer memory to approximate the inverse of the Hessian B . So instead of storing a dense $n \times n$ matrix, $L - BFGS$ keeps a record of the last m iterations where m is a small number that is chosen according to the problem at hand.³ It is for this reason that during the first m iterations, $BFGS$ and $L - BFGS$ produce exactly the same search directions.

Because of this construction, the $L - BFGS$ algorithm is less computationally intensive and it only requires $\mathcal{O}(mn)$ operations per iteration. So it is much better suited for problems where the number of dimensions n is large. For this reason it is the algorithm of choice in this thesis.

2.3 L-BFGS-B

Finally $L - BFGS - B$ comes naturally as an extension of $L - BFGS$. The B stands for the inclusion of Boundaries. $L - BFGS - B$ requires two extra steps on top of $L - BFGS$. First, a gradient projection that reduces the dimensionality of the problem. Depending on the problem, the gradient projection could potentially save a lot of iterations by eliminating those variables that are at bound at the optimum. After that, the subspace minimization; here, during the search step phase the step length is restricted as much as necessary in order to remain within the box defined by 1.1.

2.3.1 Gradient Projection

The original algorithm was created for the case when n is large and f is smooth. Its first step is a gradient projection similar to the one outlined in [4, 5] which is used to determine an active set corresponding to those variables that are on either their lower or upper bounds. The active set defined at point x^* is:

$$\mathcal{A}(x^*) = \{i \in \{1 \dots n\} | x_i^* = l_i \vee x_i^* = u_i\} \quad (2.2)$$

Working with this active set is efficient in large scale problems. A pure line search algorithm would have to choose a step length short enough to remain within the box defined by u_i and l_i . So if at the optimum, a large number \mathcal{B} of variables are either

³In this thesis $m < 20$, and in practice numbers between 5 and 10 are regularly used. There is no way of knowing a priori what choice of m will provide the best results

on the lower or the upper bound. At least a number \mathcal{B} of iterations might be needed. Gradient projection tries to reduce this number of iterations. In the best case, only 1 iteration is needed instead of \mathcal{B} .

Gradient projection works on the approximation model:

$$m_k(x) = f(x_k) + \nabla f(x_k)^T(x - x_k) + \frac{(x - x_k)^T H_k(x - x_k)}{2} \quad (2.3)$$

where H_k is a $L - BFGS - B$ approximation to the Hessian $\nabla^2 f$ stored in the implicit way defined by $L - BFGS$.

In this first stage of the algorithm a piecewise linear segment starts on the current point x_k in the direction $-\nabla f(x_k)$. Whenever this direction encounters one of the constraints, the line segment turns corners in order to remain feasible. The path is nothing but the feasible piecewise projection of the negative gradient direction on the constraint box determined by the values \vec{l} and \vec{u} . At the end of this stage, the value of x that minimizes $m_k(x)$ restricted to this piecewise gradient projection path is known as the 'Cauchy point' x^c .

2.3.2 Subspace Minimization

The problem with gradient projection is that its search direction does not take advantage of information provided implicitly by H_k , and therefore the speed of convergence is at best linear. It is for this reason that a stage two is necessary. Stage 2 or subspace minimization uses an $L - BFGS$ implicit approximation of the Inverse Hessian matrix restricted to the free variables that are not in the active set $\mathcal{A}(x^c)$.

The idea at a higher level is to solve the constrained problem 2.3, but only on those dimensions that are free. The starting point for this new problem will be the previously found Cauchy point x^c , the $L - BFGS$ approximation will provide a new search direction \hat{d}^u that takes implicit advantage of second order approximations of the Hessian matrix.

The algorithm will move in the direction $\hat{d}^* = \alpha^* \hat{d}^u$ where α^* (the step length) is chosen so that the new point \bar{x}_i satisfies weak wolfe descent and curvature conditions, A third restriction on the step length is added so that the next iteration stays feasible. Once this step is finished, the next and final step will be the termination condition. If the termination condition fails a new gradient projection and subspace minimization will be needed.

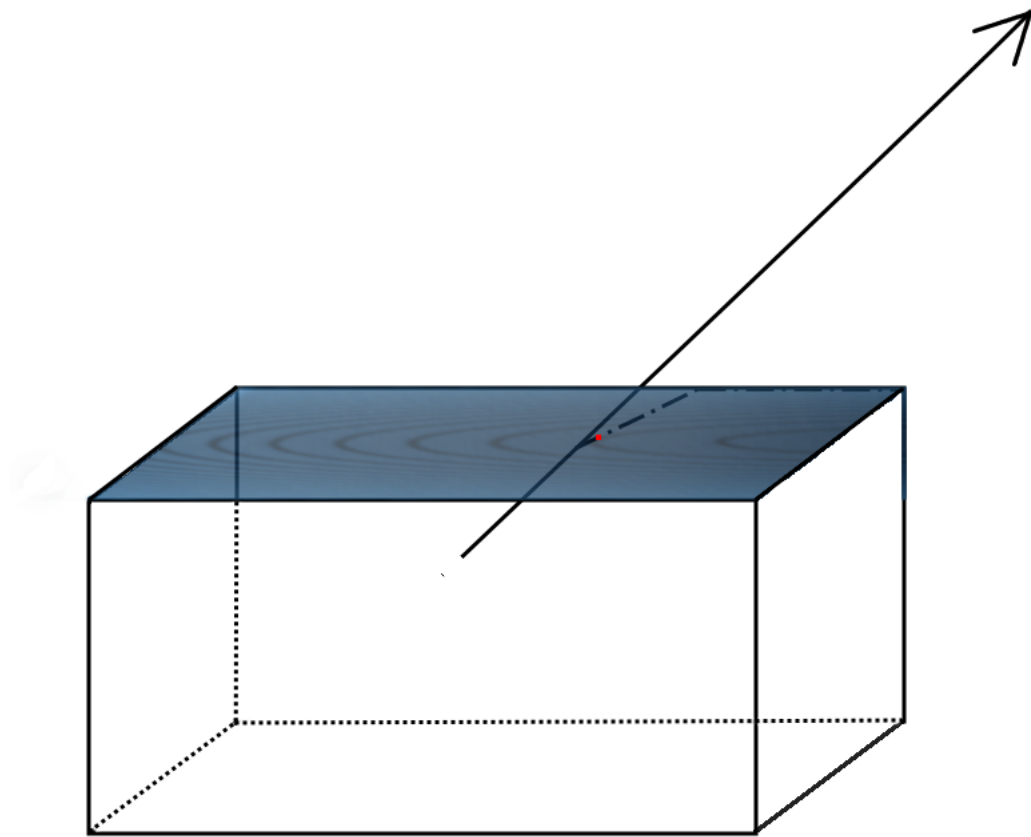


FIGURE 2.1: The arrow represents the direction of the gradient. The dotted path represents the projected gradient path from the gradient and onto the box. The region represents the level sets of the model. The optimal point (in red) is the Cauchy point x^c

Chapter 3

Modifications to the L-BFGS-B algorithm

We made three main changes to the original $L-BFGS-B$ algorithm. They concern the line search Wolfe conditions, the line search methodology, and the termination condition.

3.1 The Wolfe conditions

Probably the most important change made to the original code was the change in the curvature condition. It is accepted that the Wolfe conditions work very well whenever the function f is smooth [6]. Originally there are two Wolfe conditions: one of them is the Armijo condition, also known as the sufficient decrease requirement. The other one is the curvature condition, of which the most popular version is the “strong Wolfe” curvature condition:

$$|d_k^T \nabla f(x_k + \alpha_k d_k)| \leq |d_k^T \nabla f(x_k)| \quad (3.1)$$

Where d_k represents the search direction. The strong Wolfe condition is natural for smooth optimization. Its goal is to find a step length long enough that the slope has been reduced “sufficiently”, but the problem is that it does not work well for the Non-Smooth case. This is because near the minimal points, there may be abrupt changes in the gradient. For example the curvature condition can never be satisfied when $f(x) = |x|$, because the slope never becomes flat. The “weak wolfe” condition

$$d_k^T \nabla f(x_k + \alpha_k d_k) \geq d_k^T \nabla f(x_k) \quad (3.2)$$

is all that is needed to guarantee that the *BFGS* updated inverse Hessian approximation is positive definite [7]. This weak version is the one that will be used in this thesis. It was changed inside of the line search algorithm explained in the next section.

3.2 The line search methodology

The original *FORTRAN* code [2] contains a line search subroutine. This subroutine first establishes a maximum step. This maximum step is such that it guarantees that the step length stays within the bounding box delimited by l and u . This step maximum is defined in between lines 2607 and 2633 of the *lbfgsnomessages.f90* file. This part of the line search was not changed in this thesis. The next part of the line search is finding a step length that satisfies a decrease and a curvature condition (aka. Armijo and Wolfe conditions). The next procedure was completely changed for the purpose of this thesis. and is only explained here with the purpose of showing why it needed to be modified. It was called on line 2662 of *lbfgsnomessages.f90* and it has been left commented out in order to show how this call was before the modifications.

The original line search was a two stage process. In the first stage the interval was chosen so that it contains a minimizer of the modified function, which is nothing but a modification of the Armijo condition.

$$\Psi(\alpha_k) = f(x_k + \alpha_k d_k) - f(x_k) - c1 \alpha_k d_k^T \nabla f(x_k)$$

If $\Psi(\alpha_k) < 0$ and $\alpha_k d_k^T \nabla f(x_k) > 0$. The interval is chosen so that it contains a minimizer of f . This first stage took place in the subroutine *dcsrch*; specifically on lines 3687 and 3709

The second stage was called on line 3713 and its mission was to find the step that satisfied Armijo and “Strong” Wolfe conditions on the original function (as opposed to the modified from stage 1) f . Both stage 1 and stage 2 called the subroutine *dcstep*. The subroutine still exists on file *lbfgsbnomessages.f90* for illustration. Although it is never called in this thesis.

But there is a problem with the function *dcstep*. It turns out that function *dcstep* was designed to work only with smooth functions in mind. The code takes advantage of quadratic and cubic approximations to the function in order to calculate the next points that satisfy armijo and wolfe conditions. Unfortunately, these second and third order approximations do not work in the Non-Smooth case, and the optimizer crashes under

the line search as it was formulated. Function *dcstep* starts on line 3779 and a sample of the approximations is shown in between lines 3881 and 3902 of *lbfgsbnomessages.f90*.

The approach to this issue is to use a line search similar to the one in *hanso* [8]. The HANSO approach is to double the step length while the Armijo condition is violated (while the decreasing conditions has not been satisfied). And once the interval has been bracketed, to do a bisection until both Armijo and Wolfe conditions are satisfied. The only difference with the HANSO approach in this thesis is that the line search in HANSO can double its step length up to 30 times ¹. Whereas in this thesis, the step length can double only as long as the step length is less than the maximum value that guarantees feasibility of the solution. This version of the bisection and expansion is found in between lines 4425 and 4454 of *lbfgsbnomessages.f90*

3.3 The termination condition

One important requirement of a practical algorithm is that it ends in a finite time. For the case of smooth functions, the usual way to check whether the algorithm has converged, is by using the projected gradient which is nothing but the projection of the negative gradient onto the bounding box defined by l and u . If this projected gradient has a small norm the algorithm has converged. In the case of nonsmooth functions however, this is not necessarily true and the function at the minimum, may have a wedge. In this wedge the projected gradient may not vanish. Furthermore, if there is a sequence of points that approaches the optimum x from the right, the projected gradients corresponding to this sequence of points might be completely different from the projected gradients associated to a sequence of points that approach the optimum x from the left.

Given this set of conditions, there is the need for a special set of rules to establish the finalization of each optimization.

3.3.1 Description of the Solution

Overton and Lewis formulate an algorithm that gives a practical solution to this problem on section 6.3 [7] The best practical methodology should be to calculate wheter zero or a very small number is the norm of a vector that is part of the convex hull of gradients evaluated at points near the optimum candidate x . In order to make sure that the gradient zero $\vec{0}$ or a vector with a very small norm smaller than a very tiny

¹for all practical purposes this is considered as a good limit of iterations

tolerance $\tau_d > 0$ is part of the convex hull calculated near a neighbourhood of the optimum. The algorithm needs to keep a record of the latest gradient vectors in this small neighbourhood of the point where it is suspected that the optimum is located. The neighborhood is defined as those points with a distance to x smaller than a small tolerance $\tau_x > 0$ and no more than $J \in \mathbb{N}$ iterations back. This list of gradients is referred to as the set \mathcal{G} [7]

With this list \mathcal{G} of gradients at hand. The next step is to solve a quadratic program that locates the vector with the minimal norm that lives in the convex hull of these gradients. If the norm of this vector has a norm smaller than τ_d , the algorithm ends with a message of convergence success. If the minimum such norm is larger than the tolerance, the algorithm must continue to the following iteration and not terminate.

Every vector in the convex hull can be expressed as a nonnegative linear combination Gz of those vectors in \mathcal{G} . Where G is the matrix with columns made up of gradients in \mathcal{G} ; and z is such that $\sum_{i=1}^n z_i = 1$ and $z_i \geq 0$.

In other words the objective is to find the right combination of z that minimizes the norm $\|Gz\|_2$. This is equivalent to solving the optimization problem

$$\begin{aligned} \min \quad & q(z) = \|Gz\|_2^2 = z^T G^T G z \\ \text{s.t.} \quad & \sum_{i=1}^J z_i = 1 \\ & z_i \geq 0. \end{aligned} \tag{3.3}$$

The solution to this problem z^* has the associated vector Gz^* . And if $\|Gz^*\| < \tau_d$ the algorithm converges.

3.3.2 The Solution of the Quadratic Program

The solution of 3.3 was implemented with a practical primal-dual methodology. This methodology is the same methodology implemented by Skajaa [9] in his thesis. His code `qpspecial` was implemented in *FORTTRAN* and is part of *lbfgsbnomessages.f90*. The solution is a typical Mehrotra's Predictor-Corrector algorithm applied to quadratic programming.

This algorithm initially tries to solve the *KKT* conditions. The *KKT* is a system of equations whose solution identifies the solution of 3.3.

In order to solve this system of equations a variant of Newton's method is used. This Newton's method as usual, provides a search direction which requires a corresponding

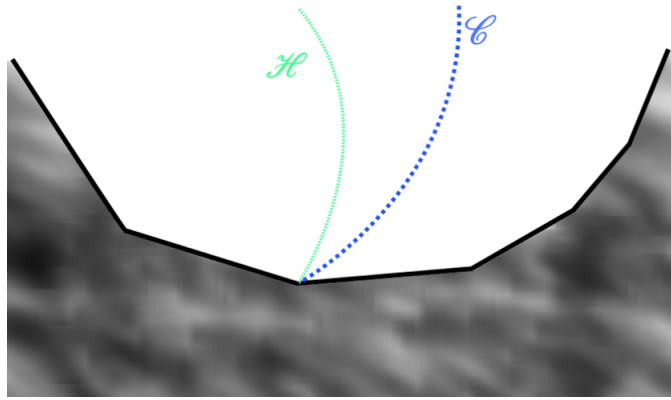


FIGURE 3.1: The central path \mathcal{C} is approached from a current and noncentral point on the tangent path \mathcal{H}

step length. Since this is an interior point method, it approaches the solution following a path inside the convex interior of the feasible set.

The problem is that the step length is usually very small because otherwise the condition of $z > 0$ or $s > 0$ where s is the slackness variable is violated. The method by itself approaches the solution very close to the boundaries of the feasible set. Not close to the *central path* where the step lengths could be larger and the convergence therefore would be faster. It is for this reason that a centering of the path is performed. This first step is known as the “predictor” step.

Besides the centering; Mehrotra proposed a “correction” step which pushes the solution closer to the central path by taking into account the curvature of the approaching *central path*. This step is also used to calibrate some of the parameters used during the centering in the “predictor” stage.

Chapter 4

the solution test functions

In order to make some tests, a few functions will be evaluated. The most important function to test this non-smooth optimizer is a modified version of rosenbrock's:

$$f(x) = (x_1 - 1)^2 + \sum_{i=2}^n |x_i - x_{i-1}^2|^p \quad (4.1)$$

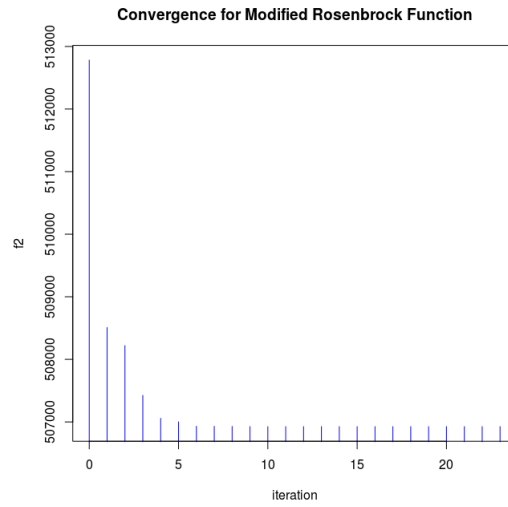
Where the value of p changes the behaviour of the optimizer. This function can be proven to be lipschitz continuous whenever $p > 1$ if restricted to the domain defined by

$$x_i = \begin{cases} [-100, 100] & \text{if } i \in \text{even numbers} \\ [10, 100] & \text{if } i \in \text{odd numbers} \end{cases} \quad (4.2)$$

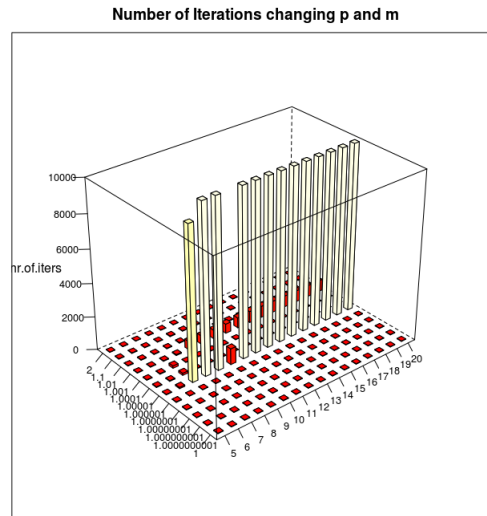
and in fact, whenever the function is restricted to a finite domain this function will be lipschitz continuous for $p > 1$. Whenever $p > 1$ the function $f(z) = |z|^p$ is zero 0 around zero because the derivative $p|z|^{p-1}$ is zero whenever z tends to zero from the right. (this is also the case from the left because it is an even function). However the second derivative will not be as nice.

For the case when $p \leq 1$ the second derivative tends to infinity. $\lim_{x \rightarrow 0^+} f' = \infty$. Which is already well known given the "heavyside" look of $f(z) = |z|$.

The convergence of the algorithm smoothly descends to the objective



The converge is adversely affected by the selection of p as one would expect. Values of p descending to 1 make the function less "smooth" and have the adverse effect of making the convergence much more difficult. In this exercise it is noticeable how slow the convergence becomes for a few specific values of p . In particular for 1.0001



$$f(x_k + \alpha_k p_k) \leq f(x_k) + c_1 \alpha_k p_k^T \nabla f(x_k) \quad (4.3)$$

Appendix A

Appendix Title Here

Write your Appendix content here.

Bibliography

- [1] Richard H. Byrd, Peihuang Lu, Jorge Nocedal, and Ci You Zhu. A limited memory algorithm for bound constrained optimization. *SIAM J. Sci. Comput.*, 16(5):1190–1208, 1995. ISSN 1064-8275. doi: 10.1137/0916069. URL <http://dx.doi.org/10.1137/0916069>.
- [2] Ciyu Zhu, Richard Byrd, Jorge Nocedal, and Jose Luis Morales. Lbfgsb 3.0. <http://www.ece.northwestern.edu/~nocedal/lbfgsb.html>, 2011.
- [3] Jorge Nocedal and Stephen J. Wright. *Numerical optimization*. Springer Series in Operations Research. Springer-Verlag, New York, 1999. ISBN 0-387-98793-2. doi: 10.1007/b98874. URL <http://dx.doi.org/10.1007/b98874>.
- [4] A. R. Conn, N. I. M. Gould, and Ph. L. Toint. Global convergence of a class of trust region algorithms for optimization with simple bounds. *SIAM J. Numer. Anal.*, 25(2):433–460, 1988. ISSN 0036-1429. doi: 10.1137/0725029. URL <http://dx.doi.org/10.1137/0725029>.
- [5] Jorge J. Moré and Gerardo Toraldo. Algorithms for bound constrained quadratic programming problems. *Numer. Math.*, 55(4):377–400, 1989. ISSN 0029-599X. doi: 10.1007/BF01396045. URL <http://dx.doi.org/10.1007/BF01396045>.
- [6] Dong-Hui Li and Masao Fukushima. On the global convergence of the BFGS method for nonconvex unconstrained optimization problems. *SIAM J. Optim.*, 11(4):1054–1064 (electronic), 2001. ISSN 1052-6234. doi: 10.1137/S1052623499354242. URL <http://dx.doi.org/10.1137/S1052623499354242>.
- [7] Adrian S. Lewis and Michael L. Overton. Nonsmooth optimization via quasi-Newton methods. *Math. Program.*, 141(1-2, Ser. A):135–163, 2013. ISSN 0025-5610. doi: 10.1007/s10107-012-0514-2. URL <http://dx.doi.org/10.1007/s10107-012-0514-2>.
- [8] Michael Overton. Hanso 2.02. <http://www.cs.nyu.edu/faculty/overton/software/hanso/>, 2012.

-
- [9] Anders Skaaja. Limited memory bfgs for nonsmooth optimization. Master's thesis, New York University, Courant Institute of Mathematical Sciences, 251 Mercer Street, New York, NY 10012, 2010.