

# Implementation of High Precision Arithmetic in the BFGS Method for Nonsmooth Optimization

Allan Kaku

New York University, Courant Institute of Mathematical Sciences

251 Mercer Street

New York, NY 10012

January 2011

A thesis submitted in partial fulfillment  
of the requirements for the degree of  
master's of science  
Department of Mathematics  
New York University  
January 2011

-----  
Advisor: Michael L. Overton

### **Abstract**

The results of this thesis investigate the claim that the effectiveness of the BFGS variable metric (quasi-Newton) method with an inexact weak Wolfe line search for minimizing nonsmooth, nonconvex functions is limited by machine precision, rather than defects in the algorithm itself. To investigate this question, the algorithm was implemented using double-double precision, i.e. twice as precise as floating-point double. The double and double-double precision implementations were tested on Nesterov's Chebyshev-Rosenbrock functions, which are known to be difficult to solve. Results for these and other test functions for both the normal and high-precision BFGS codes are presented.

## Contents

<b>1</b>	<b>Introduction</b>	<b>6</b>
<b>2</b>	<b>BFGS</b>	<b>6</b>
<b>3</b>	<b>High Precision Arithmetic</b>	<b>7</b>
3.1	Double-double . . . . .	8
<b>4</b>	<b>Implementation</b>	<b>9</b>
4.1	MATLAB to C . . . . .	9
4.2	Elimination of CBLAS and CLAPACK . . . . .	10
4.3	C to C++ and overload functions . . . . .	11
<b>5</b>	<b>Test Functions</b>	<b>11</b>
5.1	F Functions . . . . .	11
5.2	T Functions . . . . .	13
5.3	Nesterov's Chebyshev-Rosenbrock Functions . . . . .	14
<b>6</b>	<b>Results and Comparison of Performance</b>	<b>15</b>
6.1	F Functions Results . . . . .	16
6.2	T Functions Results . . . . .	17
6.3	Nesterov Functions Results . . . . .	19
<b>7</b>	<b>Conclusion</b>	<b>25</b>
	<b>Acknowledgements</b>	<b>27</b>
	<b>References</b>	<b>28</b>

## List of Figures

3.1	<i>Algorithm for double-double addition. . . . .</i>	9
5.1	<i>Contour plot for NCR-NS1 (left) and NCR-NS2 (right), <math>n=2</math>. The line segments show the iterates of BFGS initialized from 7 random starting points. . . . .</i>	15
6.1	<i>T3 performance plots for <math>n = 10</math> (left) and <math>n = 50</math> (right), with the final DD function values sorted and plotted with the corresponding final D function value for 1000 random starting points. It appears that there are multiple local minima. . . . .</i>	18
6.2	<i>T4 performance plots for <math>n = 10</math> (left) and <math>n = 50</math> (right), with the final DD function values sorted and plotted with the corresponding final D function value for 1000 random starting points. It appears that there are multiple local minima. . . . .</i>	19
6.3	<i>NCR-NS1 performance plots, with the final DD function values sorted and plotted with the corresponding final D function value for 1000 starting points. Top left: Plot for <math>n = 2</math>. The DD code reduces the function value to as low as its precision allows. Top right: Plot for <math>n = 3</math>. The DD code overall does better than the D code, but with a much larger range. Middle left: Plot for <math>n = 4</math>. The DD code again does better than the D code, but final values are much closer. Middle right: Plot for <math>n = 5</math>. Both codes do equally well in some runs, and both codes have a large range of final values. Bottom left: Plot for <math>n = 6</math>. A limit is placed on the maximum number of iterations for <math>n \geq 6</math>. Neither code does very well, with little difference between the two. Bottom right: Plot for <math>n = 7</math>. Neither code does very well, with little difference between the two. . . . .</i>	21

- 6.4 *NCR-NS2 performance plots, with the final DD function values sorted and plotted with the corresponding final D function value for 1000 starting points.*  
*Top left: Plot for  $n = 2$ . The DD code finds the minimizer much more often than the D code. Top right: Plot for  $n = 3$ . The DD code finds the minimizer most of the time, while the D code converges to the other Clarke stationary points more often. Top-Middle left: Plot for  $n = 4$ . The DD code finds the minimizer rather than the Clarke stationary points much more often than the D code. Top-Middle right: Plot for  $n = 5$ . The DD code finds a lower stationary point than the D code almost every time, but only finds the minimizer on a quarter of the runs. Bottom-Middle left: Plot for  $n = 6$ . Neither code finds the minimizer very often, but the DD code still finds a lower stationary point most of the time. Bottom-Middle right: Plot for  $n = 7$ . Neither code does very well in finding the minimizer a majority of the time, with minimal improvement of the DD code over the D code. Bottom: Plot for  $n = 8$ . Both codes perform practically the same. . . . . 22*
- 6.5 *Function values after each iteration of several runs of NCR-NS2.*  
*Top left: Plot for  $n = 3$ . The DD code is able to find a smaller stationary point than the D code after several more iterations beyond the D code termination. Top right: Plot for  $n = 5$ . The DD code is able to get past several more corners than the D code. Bottom: Plot for  $n = 7$ . The DD code gets past many more corners than the D code. 25*

## List of Tables

- |   |  |    |
|---|--|----|
| 1 | <i>Minimum and maximum final values found by the D and DD codes for the F functions, <math>n = 10</math>, and the optimal value. . . . .</i>                       | 16 |
| 2 | <i>Minimum and maximum final values found by the D and DD codes for the F functions, <math>n = 50</math>, and the optimal value. . . . .</i>                       | 16 |
| 3 | <i>Minimum and maximum final values found by the D and DD codes for the F functions, <math>n = 200</math>, and the optimal value. . . . .</i>                      | 16 |
| 4 | <i>Minimum and maximum final values found by the D and DD codes for the T functions, <math>n = 10</math>, and the optimal value. . . . .</i>                       | 17 |
| 5 | <i>Minimum and maximum final values found by the D and DD codes for the T functions, <math>n = 50</math>, and the optimal value. . . . .</i>                       | 17 |
| 6 | <i>Minimum and maximum final values found by the D and DD codes for the T functions, <math>n = 200</math>, and the optimal value. . . . .</i>                      | 18 |
| 7 | <i>Values of BFGS for double and double-double precision. Superlinear convergence can be seen in the higher iterations, before the program terminates. . . . .</i> | 20 |

## 1 Introduction

The applicability of the Broyden-Fletcher-Goldfarb-Shanno (BFGS) method [NW06] is in solving the unconstrained optimization problem

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

BFGS is an iterative method that generates a sequence of points to find a local minimizer  $x^*$ , which may not be a global minimizer due to nonconvexity of  $f$ . In this thesis, the function  $f$  is only assumed to be continuous, and not continuously differentiable. For an introduction to algorithms for solving nonconvex, nonsmooth optimization problems, see [Kiw85].

In work by Lewis and Overton [LO10] and Skajaa [Ska10], BFGS and the limited memory variation LBFGS were shown to work well in solving nonsmooth test problems, with the primary limitation apparently being machine precision. By implementing a high-precision version of the BFGS solver, this thesis investigates the claim that the algorithms are unable to solve difficult problems because of rounding errors, rather than because of a breakdown in the algorithm itself.

Section 2 will give a brief discussion of the BFGS method. In section 3, a description of double-double high-precision arithmetic is provided, as developed by David Bailey in [Bai10], and distributed by Bailey and Xiaoye (Sherry) Li. Section 4 details the implementation process of incorporating double-double into the BFGS code, using the C code written by Skajaa [Ska10], explaining the changes made and their implications for the results. Section 5 describes the test functions that were used to compare the double and double-double versions, including Nesterov's Chebyshev-Rosenbrock functions. Finally, section 6 gives the results of the tests, comparing the correctness and accuracy of the outputs, as well as the time taken.

A website<sup>1</sup> is being developed with the BFGS code used in this thesis, which includes both the double precision version written by Anders Skajaa and the new double-double version of the code. The double-double version of BFGS uses C++, in order to allow the use of classes with function and operator overloading. In addition, the test problems used to verify the algorithms and the Nesterov functions are available in both double and double-double precision.

## 2 BFGS

BFGS is the most popular quasi-Newton update formula. Specifically, we will focus on the class of line-search methods that use the search direction defined by

$$d^{(j)} = -C_j \nabla f(x^{(j)})$$

---

<sup>1</sup><http://cims.nyu.edu/~aak357/>

where  $C_j$ , an approximation to the inverse Hessian  $\nabla^2 f(x^j)$ , is updated by a quasi-Newton update formula with certain properties. The BFGS update formula is

$$C_{j+1} = (I - p_j s^{(j)} (y^{(j)})^T) C_j (I - p_j y^{(j)} (s^{(j)})^T) + p_j s^{(j)} (s^{(j)})^T$$

where

$$p_j = ((y^{(j)})^T s^{(j)})^{-1}$$

and

$$\begin{aligned} s^{(j)} &= x^{(j+1)} - x^{(j)} \\ y^{(j)} &= \nabla f(x^{(j+1)}) - \nabla f(x^{(j)}). \end{aligned}$$

When  $f$  is continuously differentiable, the BFGS method exhibits superlinear convergence after enough iterations, under standard regularity assumptions. See Nocedal and Wright [NW06] for more details.

The Limited-Memory BFGS method, or LBFGS, is less computationally intensive, requiring  $\mathcal{O}(mn)$  floating point operations per update rather than the  $\mathcal{O}(n^2)$  floating point operations required for BFGS, where for the first  $m$  iterations, LBFGS and BFGS generate the same search directions. The best choice of  $m$  is problem-dependent, and the method may require more iterations to converge than BFGS. Tests using LBFGS were not completed in this thesis; however similar experimentation is planned for the future.

In dealing with the nonsmooth case, BFGS typically succeeds in finding a local minimizer, despite being developed for smooth functions. However, the right line search must be used, namely, a weak Wolfe line search [LO10, Ska10]. The rate of convergence in the nonsmooth case is typically linear, not superlinear.

In our tests, the methods begin with a random starting point, unless otherwise specified. A set of points are generated using the `rand` and overloaded `rand` function in C/C++ and normalizing the output between  $[-1, 1]$ . Thus, the starting point  $x^{(0)} \in \mathbb{R}^n$  is drawn randomly from the uniform distribution on  $[-1, 1]^n$ .

### 3 High Precision Arithmetic

High precision computing has existed for many years, and as such there are several implementations and variations in different computing languages. In MATLAB, for instance, the Symbolic Math Toolbox contains a function called Variable Precision Arithmetic, or `vpa`, which allows users to specify the number of digits of accuracy they desire<sup>2</sup>. MATLAB automatically ensures that the results are precise to at least the number specified. While this allows high precision computing, it is slow due to the method by which the extra precision computations are performed, as well as the general slower nature of MATLAB in executing code.

<sup>2</sup><http://www.mathworks.com/help/toolbox/symbolic/vpa.html>

In other programming languages, many high precision libraries have been created. They can be split into two distinct groups, based on how the precision is represented. One group uses a multiple-digit format, which stores a sequence of digits or bits with a single exponent. The group we will focus on is a multiple-component format, which stores values as an unevaluated sum of normal precision numbers. The advantage of this second representation is that calculations are done much faster. Computations are done on a component level, and then combined at the end to gain the extra precision. Thus, calculations are on the order of magnitude of normal machine precision. This idea was apparently originally suggested by Priest [Pri91]. Two examples of this multiple-component format are double-double and quad-double precision floating point arithmetic, implemented by Hida, Li, and Bailey [LHB01]. As fixed precision libraries, both run faster than other arbitrary precision libraries. We will be focusing on the double-double library.

### 3.1 Double-double

A double-double number is the sum of two IEEE double format numbers [Ove01]. The double-double number  $(a_0, a_1)$  represents the exact sum  $a = a_0 + a_1$ . The first number,  $a_0$ , is the most significant component, with  $a_1$  providing the extra precision. Since a 64 bit IEEE double format number uses 1 bit for the sign, 11 bits for the exponent and 52 bits for the significand, and since there is an additional "hidden bit" for normalized numbers as explained in [Ove01], it effectively has 53 bits of precision, or approximately 16 decimal digits. An example of a double-double number would be  $\pi + e * 10^{-16}$ , namely  $3.14159265358979351029082622918449e+00$ , which has 106 bits of precision or approximately 32 digits. In fact, because of the binary nature of the representation, the number  $1 + 2^{-1000}$  can be stored *exactly* as a double-double number, since both the numbers 1 and  $2^{-1000}$  are exact floating point double format numbers. Although one might say that such a number has more than 106 bits of precision, this is somewhat artificial as most of the internal zero bits are not stored, and if arithmetic were done using such numbers, the least significant bits would be lost rapidly. Thus, double-double numbers effectively have 106 bits of precision, or about 32 digits, somewhat less than if they were stored using a 128-bit floating point word, but the advantage is that arithmetic on double-doubles can be done rapidly using standard double-format hardware.

This representation was implemented in C++ by creating a new class, overloading the arithmetic and comparative operators and standard mathematical functions [LHB01]. In addition, several constants and random number generators were implemented. The source code and documentation is all available open source at <http://crd.lbl.gov/~dhbailey/mpdist/>. The necessary libraries and header files to use double-double are already packaged with BFGS in the high-precision version of the code on the companion website to this thesis.

The algorithms to execute the operations between double-doubles is complex, and for a complete explanation, consult [LDB<sup>+</sup>02]. The double-double addition algorithm



```

void ddadd(double dda[2], double ddb[2], double ddc[2])
{
    /* Compute double-double = double-double + double-double
       (ddc = dda + ddb). */
    /* dda[0] represents high-order word, dda[1] represents
       low-order word. */
    double bv;
    double s1, s2, t1, t2;

    /* Add two high-order words. */
    s1 = dda[0] + ddb[0];
    bv = s1 - dda[0];
    s2 = ((ddb[0] - bv) + (dda[0] - (s1 - bv)));

    /* Add two low-order words. */
    t1 = dda[1] + ddb[1];
    bv = t1 - dda[1];
    t2 = ((ddb[1] - bv) + (dda[1] - (t1 - bv)));

    s2 += t1;

    /* Renormalize (s1, s2) to (t1, s2) */
    t1 = s1 + s2;
    s2 = s2 - (t1 - s1);

    t2 += s2;

    /* Renormalize (t1, t2) */
    ddc[0] = t1 + t2;
    ddc[1] = t2 - (ddc[0] - t1);
}

```

Figure 3.1: *Algorithm for double-double addition.*

shown in Figure 3.1 requires 20 floating-point operations, and therefore runs about an order of magnitude slower than double arithmetic.

## 4 Implementation

We began with a MATLAB code for BFGS developed M. Overton as a part of the HANSO (Hybrid Algorithm for Nonsmooth Optimization) package<sup>3</sup>. The process of incorporating the double-double class into the source code for BFGS went through several iterations. These alterations to the original MATLAB code are detailed below, as well as their impact on the algorithm’s behavior and results.

### 4.1 MATLAB to C

The first translation to the original MATLAB code to C was done by A. Skajaa. This initial translation of the HANSO package into C resulted in a much faster running version, although it was still partly dependent on MATLAB to run. While the BFGS code and its required functions were translated to C, the HANSO code still ran through MATLAB, using MEX as an interface to call the C code. Furthermore, the test functions were also called through MATLAB. In order to properly implement the

<sup>3</sup><http://www.cs.nyu.edu/faculty/overton/software/hanso/>

double-double C++ library, a stand-alone version BFGS, and eventually HANSO, would need to be created.

Taking just the BFGS source code, the program was stripped down to only its essential components. This included implementing a function pointer to the various test functions, now translated into C. To simplify the program and make it more portable, it was necessary to download and dynamically link the necessary libraries and header files from CLAPACK<sup>4</sup> and CBLAS<sup>5</sup> that are required to do the matrix and vector operations. An executable was also created in order to verify the correct translation of the test functions and alterations made to the source code, as well perform the tests for this thesis. A similar translation can be, but has not been yet, implemented in the Gradient Sampling code, which is another component of HANSO.

In addition, the test functions were translated from MATLAB to C, with the T functions translated by A. Skajaa, and the F functions and Nesterov functions translated by the author, who also took responsibility for the debugging of all the code.

## 4.2 Elimination of CBLAS and CLAPACK

Although we now had a version of BFGS in C, the double-double library could not yet be added. A majority of the major calculations within BFGS were performed, for efficiency and speed, using the optimized algebraic operations provided by CBLAS and CLAPACK. While useful in both its ease of implementation and its recognition as a standard in numerical computing, it took away the control over the calculations, which was necessary in adding double-double. Although Li has been developing XBLAS, an Extended and Mixed precision BLAS standard [LDB<sup>+</sup>02], it does not yet include all the necessary operations for BFGS. Thus, in order to fully implement double-double precision in all the calculations, the function calls to CBLAS and CLAPACK were replaced with hard-coded operations that executed the necessary calculations. While resulting in slower runtimes, as these codes were not at all optimized, they paved the way for the use of double-double precision. In addition, these changes simplified the program and made it more portable by eliminating the need for the CBLAS and CLAPACK libraries and header files.

In addition, the QP stopping criterion in BFGS [LO10] was removed. The main purpose of the criterion was to stop the code if the solution was only minimally improved with further iterations, thus reducing the running time of the program. However, with the anticipation of using higher precision, allowing slight improvements on each iteration would be required in order for the extra precision to affect the output. Therefore, the QP stopping criterion was commented out of the code in order to test the increased accuracy that double-double provides. It is recommended it remains excluded when running any functions using double-double precision.

---

<sup>4</sup><http://www.netlib.org/clapack/>

<sup>5</sup><http://www.netlib.org/blas/>

### 4.3 C to C++ and overload functions

While C and C++ code are for the most part completely compatible with one another, they differ most importantly in the functionality of classes in C++. Li's double-double library is primarily written in C++ using a special class called `dd_real` to represent double-double numbers, but also contains low-level C implementations of the operations, without the use of classes. Rather than creating a new class, the C version stores each number in an array of size two, which represents the double-double number, as in Figure 3.1 on page 9. Implementing this complicated method of arrays and low-order functions, while possible, made the code unnecessarily long and convoluted. A simpler approach was to first convert the BFGS C source code to C++. A few changes were made in order for this to be possible, namely substituting the use of the `new` operator in lieu of `malloc`, as well as the changing of a few header files and print statements. In general, the source code was not changed in any significant way during the translation.

Once compiled in C++, the double-double libraries<sup>6</sup> were linked and the necessary header files added. Because of the use of classes and overloaded functions, changing the precision of the necessary variables was as simple as declaring them of type `dd_real` instead of `double`. Once a few initial tests were run on simple functions, and the program was verified as operating correctly, more comprehensive tests were made on less trivial functions. All tests were performed using the same random starting points for the double and double-double versions in order to create meaningful comparisons.

## 5 Test Functions

In order to compare the performance of the double (D) precision code to the double-double (DD) precision code, the algorithms were implemented on several test functions.

### 5.1 F Functions

The first set of test functions are nine problems from [HMM04], denoted F1-F9, and are all nonsmooth at the minimizer. The first five F1-F5 are convex and the other four F6-F9 are nonconvex. They are defined as follows.

**F1:** Generalization of MAXQ

$$f(x) = \max_{1 \leq i \leq n} x_i^2$$

---

<sup>6</sup>Actually, the quad-double libraries were used, which includes the double-double class.

**F2:** Generalization of MAXHILB

$$f(x) = \max_{1 \leq i \leq n} \left| \sum_{j=1}^n \frac{x_j}{i+j-1} \right|$$

**F3:** Chained LQ

$$f(x) = \sum_{i=1}^{n-1} \max \{ -x_i - x_{i+1}, -x_i - x_{i+1} + (x_i^2 + x_{i+1}^2 - 1) \}$$

**F4:** Chained CB3 I

$$f(x) = \sum_{i=1}^{n-1} \max \{ x_i^4 + x_{i+1}^2, (2 - x_i)^2 + (2 - x_{i+1})^2, 2e^{-x_i + x_{i+1}} \}$$

**F5:** Chained CB3 II

$$f(x) = \max \left\{ \sum_{i=1}^{n-1} (x_i^4 + x_{i+1}^2), \sum_{i=1}^{n-1} ((2 - x_i)^2 + (2 - x_{i+1})^2), \sum_{i=1}^{n-1} (2e^{-x_i + x_{i+1}}) \right\}$$

**F6:** Number of Active Faces

$$f(x) = \max_{1 \leq i \leq n} \left\{ g \left( - \sum_{j=1}^n x_j \right), g(x_i) \right\}, \quad \text{where } g(y) = \ln(|y| + 1)$$

**F7:** Nonsmooth Generalization of Brown Function 2

$$f(x) = \sum_{i=1}^{n-1} \left( |x_i|^{x_{i+1}^2+1} + |x_{i+1}|^{x_i^2+1} \right)$$

**F8:** Chained Miffin 2

$$f(x) = \sum_{i=1}^{n-1} \left( -x_i + 2(x_i^2 + x_{i+1}^2 - 1) + \frac{7}{4} |x_i^2 + x_{i+1}^2 - 1| \right)$$

**F9:** Chained Crescent I

$$f(x) = \max \left\{ \sum_{i=1}^{n-1} \left( x_i^2 + (x_{i+1} - 1)^2 + x_{i+1} - 1 \right), \sum_{i=1}^{n-1} \left( -x_i^2 - (x_{i+1} - 1)^2 + x_{i+1} + 1 \right) \right\}$$

Through experiments, it has been determined that, apparently, none of these problems have multiple local minima for  $n \leq 200$  [Ska10]. All functions were tested with  $n = 10, 50, 200$  variables using random starting points for 10 runs on each code.

## 5.2 T Functions

The second set of test problems were five functions taken from [LTS<sup>+</sup>02], denoted T1-T5. All are again nonsmooth and can be defined with any number of variables. The first two, T1 and T2, are convex. The problems are defined as follows.

**T1:** Problem 2 from TEST29

$$f(x) = \max_{1 \leq i \leq n} |x_i|$$

**T2:** Problem 5 from TEST29

$$f(x) = \sum_{i=1}^n \left| \sum_{j=1}^n \frac{x_j}{i+j-1} \right|$$

**T3:** Problem 6 from TEST29

$$f(x) = \max_{1 \leq i \leq n} |(3 - 2x_i) x_i + 1 - x_{i-1} - x_{i+1}|, \quad \text{with } x_0 = x_{n+1} = 0$$

**T4:** Problem 11 from TEST29

$$f(x) = \sum_{k=1}^{2(n-1)} |f_k(x)|$$

where

$$f_k(x) = x_i + x_{i+1} ((1 + x_{i+1}) x_{i+1} - 14) - 29 \quad \text{if } \text{mod}(k, 2) = 0$$

$$f_k(x) = x_i + x_{i+1} ((5 - x_{i+1}) x_{i+1} - 2) - 13 \quad \text{if } \text{mod}(k, 2) = 1$$

$$i = \lfloor (k+1)/2 \rfloor$$

**T5:** Problem 22 from TEST29

$$f(x) = \max_{1 \leq i \leq n} \left| 2x_i + \frac{1}{2(n+1)^2} \left( x_i + \frac{i}{n+1} + 1 \right)^3 - x_{i-1} - x_{i+1} \right|$$

All functions were tested for  $n = 10, 50, 200$  variables using random starting points for 10 runs on each code.

### 5.3 Nesterov's Chebyshev-Rosenbrock Functions

The three Nesterov Chebyshev-Rosenbrock (NCR) functions are all very difficult to solve. For reasons that are only partially understood, the functions are exponentially difficult in  $n$ . This makes these the most challenging test functions for our double-double BFGS solver.

The smooth function, which we will denote as NCR-S, is:

$$\tilde{f}(x) = \frac{1}{4} (x_1 - 1)^2 + \sum_{i=1}^{n-1} (x_{i+1} - 2x_i^2 + 1)^2.$$

The first nonsmooth variation, which we will denote as NCR-NS1, is:

$$\hat{f}(x) = \frac{1}{4} (x_1 - 1)^2 + \sum_{i=1}^{n-1} |x_{i+1} - 2x_i^2 + 1|.$$

The second nonsmooth variation, which we will denote as NCR-NS2, is:

$$f(x) = \frac{1}{4} |x_1 - 1| + \sum_{i=1}^{n-1} |x_{i+1} - 2|x_i| + 1|.$$

In all three cases, the only local minimizer is  $\bar{x} = [1, 1, 1, \dots, 1]^T$ .

For the smooth variant, the starting point used is  $\hat{x} = [-1, 1, 1, \dots, 1]^T$ . This point is contained in the manifold

$$M = \{x : x_{i+1} - 2x_i^2 + 1 = 0, \quad i = 1, \dots, n-1\}$$

which also contains the minimizer  $\bar{x}$ . Because  $\tilde{f}$  is the sum of a quadratic term and a nonnegative sum whose zero set is the manifold  $M$ , minimizing  $\tilde{f}$  is equivalent to minimizing its first quadratic term on  $M$ , and it is easy to see that the only stationary point of  $\tilde{f}$  is the global minimizer  $\bar{x}$ . Due to the manifold's highly oscillatory nature, BFGS requires many iterations to converge to the minimizer.

The first nonsmooth variation is nondifferentiable at points in  $M$ , including the point  $\hat{x}$ , but as explained in [GO10], it is "partly smooth" with respect to the manifold  $M$

in the sense of [Lew03]. See Figure 5.1 (left) for a contour plot for  $n = 2$ , where the line segments show the iterates of BFGS from 7 random starting points. Minimizing this function is similarly equivalent to minimizing its first quadratic term on  $M$ , and again it has only one stationary point, namely the minimizer  $\bar{x}$ . The starting point  $x$  is chosen randomly from the uniform distribution between  $-1$  and  $1$ .

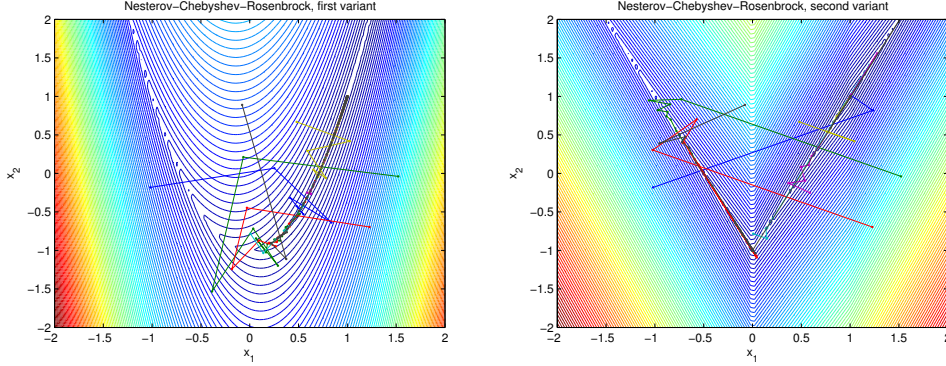


Figure 5.1: *Contour plot for NCR-NS1 (left) and NCR-NS2 (right),  $n=2$ . The line segments show the iterates of BFGS initialized from 7 random starting points.*

The second variation is not only nonsmooth, but also non-regular in the sense of Clarke (see [GO10]). Minimizing  $f$  is equivalent to minimizing its first term on the set  $S$ , where

$$S = \{x : x_{i+1} - 2|x_i| + 1 = 0, \quad i = 1, \dots, n-1\}.$$

The set  $S$  is again highly oscillatory, and contains “corners”, resulting in its non-regularity. See Figure 5.1 (right) for a contour plot for  $n = 2$ , again showing the iterates of BFGS initialized from 7 random starting points. In fact, for a given  $n$ , the NCR-NS2 function has  $2^{n-1}$  Clarke stationary points, of which all except the global minimizer  $\bar{x}$  are “corners” of  $S$  [GO10]. The starting point  $x$  is again chosen randomly from the uniform distribution between  $-1$  and  $1$ .

In all cases, the BFGS inverse Hessian approximation matrix was initialized as the identity matrix. NCR-S was tested with  $n = 2, \dots, 10$  variables, NCR-NS1 was tested with  $n = 2, \dots, 7$  variables, and NCR-NS2 was tested with  $n = 2, \dots, 8$  variables. BFGS was initialized at  $\hat{x}$  for NCR-S and at randomly generated starting points for the 1000 runs of NCR-NS1 and NCR-NS2.

## 6 Results and Comparison of Performance

To fully test the capabilities of the BFGS algorithm with higher precision, we terminate the program only when it breaks down because of rounding errors due to

machine precision. That is, all parameters are set such that the termination only occurs when a search direction is not a descent direction ( $\nabla f(x^{(j)})^T p^{(j)} < 0$ ) or a reduction is not obtained in the line search (the limit on the number of bisection or expansion steps, 100, is exceeded). For F1-F4 and F8,  $n = 200$ , the limit on the number of iterations was set to  $10^5$ , and for NCR-NS1,  $n = 6, 7$ , the limit on the number of iterations was set to  $10^{n-1}$ , so that the tests would run in a reasonable amount of time. In all other tests, the iteration limit was set large enough that it did not cause termination.

All tests were done using a dual-core 64-bit Intel Pentium Processor SU4100, 1.3 GHz, 2 MB cache, and 3 GB DDR3 memory, running Ubuntu 10.04.

## 6.1 F Functions Results

The minimum and maximum final values found by the D and DD codes over the 10 randomly generated starting points on the F functions are shown in Tables 1, 2, and 3 for  $n = 10, 50$ , and 200, respectively. Also shown is the optimal value  $f^*$  to double-double precision, when it is known.

Problem #	$f^*$	D		DD	
		min	max	min	max
F1	0	2.0066500862662096e-92	1.7598318660798305e-75	1.89733574848909154732925772036497e-231	2.29653194203331775377210255634163e-187
F2	0	7.4592726093114514e-10	6.6591901898665590e-08	9.36011982957647853299213885014951e-16	2.41962019380626166470085465814175e-14
F3	-1.272792206135785439215198517887e+01	-1.2727922061357845e+01	-1.272792206137899e+01	-1.2727922061357854392151985178865e+01	-1.2727922061357854392151985178873e+01
F4	18	1.8000000000000007e-01	1.8000000000000029e-01	1.80000000000000000003e-01	1.80000000000000000000021e-01
F5	18	1.8000000000000000e-01	1.8000000000000007e-01	1.8000000000000000000000e-01	1.80000000000000000000001e-01
F6	0	1.3322676295501871e-15	1.1324274851176532e-14	1.5942199684643182922547129492694e-32	1.36747052287869415053000131594318e-31
F7	0	2.0264389251818834e-15	9.0434354618988751e-15	6.67782176277300765732658343900183e-31	7.031855603952720389762047070724862e-31
F8	-	-6.5146142106623506e+00	-6.5146142106776281e+00	-6.5146142106776417143379112908743e+00	-6.5146142106776417143379113013307e+00
F9	0	-1.1102230246251565e-16	4.4408920985006262e-16	-3.7728666683992139855764567896852e-33	1.14446882312442438039841289255342e-31

Table 1: Minimum and maximum final values found by the D and DD codes for the F functions,  $n = 10$ , and the optimal value.

Problem #	$f^*$	D		DD	
		min	max	min	max
F1	0	4.1109297512990274e-78	5.7156919175156493e-72	4.16193040390410660681449811504549e-206	1.44435573419311235441498997967932e-180
F2	0	9.5770040422354930e-01	2.355088200918547e-07	7.32161863907410688793970137594072e-15	2.559290629588893452739575339396e-13
F3	-6.929646456281657391282747486275e+01	-6.929646456281350e+01	-6.929646456281392e+01	-6.9296464562816573912827474862540e+01	-6.929646456281657391282747486276e+01
F4	98	9.800000000000098e-01	9.8000000000000512e-01	9.8000000000000000000114e-01	9.800000000000000000000013e-01
F5	98	9.799999999999998e-01	9.8000000000000514e-01	9.799999999999999999999999999e-01	9.80000000000000000000000154e-01
F6	0	1.1546319456101562e-14	1.9628743073370842e-13	2.3671260788658168521416821333198e-31	1.53872561001547266917713004254936e-30
F7	0	2.4364638592174621e-14	5.403796513867656e-14	2.708920887289529358670188641988e-30	7.710907436775246890713125377296e-30
F8	-	-3.4795181409482488e-01	-3.4795181409546011e-01	-3.4795181409547642983703575994992e-01	-3.47951814095476429837035776089150e-01
F9	0	0	2.8865798640254070e-15	3.85185088877447170611195588516987e-33	1.193113606547892601096817833543137e-31

Table 2: Minimum and maximum final values found by the D and DD codes for the F functions,  $n = 50$ , and the optimal value.

Problem #	$f^*$	D		DD	
		min	max	min	max
F1	0	5.2745683191926849e-43	3.8387659508038081e-40	3.77374919674166623862341319931020e-43	3.6726547391480082201240313034569e-40
F2	0	1.4803212635729203e-09	9.7432916068129066e-08	2.1815668036739972169576944811633e-14	4.08757881288519634403693013811776e-13
F3	-2.814284989122491471153605611773e+02	-2.8142849891224233e+02	-2.8142849891224501e+02	-2.81428498912245914711536056117755e+02	-2.8142849891224914711536056117599e+02
F4	398	3.9800000000000233e-02	3.9800000000000495e-02	3.98000000000000000000993e-02	3.980000000000000000000012e-02
F5	398	3.9800000000000000e-02	3.9800000000000398e-02	3.98000000000000000000001e-02	3.980000000000000000000024e-02
F6	0	6.6391336872582153e-14	1.5241141682042535e-12	8.54982823487414634140425891229906e-31	6.56642127726881695146388049678088e-30
F7	0	2.4244868800304813e-13	3.281585715532334e-13	2.84629823120409264001809189794949e-29	2.19701603466485119167094768170693e-01
F8	-	-1.4086070716519524e+02	-1.4086070717278730e+02	-1.40860707172869442040965687302156e+02	-1.40860707172869442040965733244829e+02
F9	0	0	6.4392935458259079e-15	-6.62700456801495436812480410181332e-32	1.5457513733651954956727789671866e-30

Table 3: Minimum and maximum final values found by the D and DD codes for the F functions,  $n = 200$ , and the optimal value.



For most of the F functions, the D code of BFGS was able to find the minimizer fairly well when  $n = 10$ , as demonstrated in Skajaa [Ska10]. The DD code did converge to a lower final function value, reflecting the different machine precision, about  $10^{-16}$  for D and  $10^{-32}$  for DD. The exception was for F2, for which the D code only reduced the function value to  $10^{-10}$ , while the DD code was able to reduce the function value to  $10^{-16}$ . The easiest problem for either code to solve was F1, as seen by the final function values being as low as  $10^{-92}$  and  $10^{-231}$  for the D and DD code, respectively.

When  $n = 50$ , both codes exhibited the same behavior as  $n = 10$ . For  $n = 200$ , the problems were more difficult for either code to solve. On F1, both codes reduced the final function value to around  $10^{-43}$ , significantly higher than with fewer variables. Most interesting was the DD code's final function values on F7, ranging from  $10^{-29}$  to 22, while the D code always converged around  $10^{-13}$ . The reason why the DD code sometimes fails in solving F7 for  $n = 200$  needs further investigation.

## 6.2 T Functions Results

The minimum and maximum final values found by the D and DD codes over the 10 randomly generated starting points on the T functions are shown in Tables 4 and 5 on this page, and 6 on the following page for  $n = 10, 50$ , and 200, respectively. Also shown is the optimal value  $f^*$  to double-double precision, when it is known.

Problem #	$f^*$	D		DD	
		min	max	min	max
T1	0	1.4600658163256693e-16	8.2340403956125942e-16	1.73192903530815863168723237441175e-32	7.97966255429890611008730620634931e-32
T2	0	6.7428593524682112e-10	1.2019196891617859e-07	2.46273384860065925125674845594735e-15	2.08959152325082623512778323565036e-13
T3	0	3.2196467714129540e-15	4.4617856727674143e-01	2.83496887813801117569839953148502e-31	4.46176671004398955088565094513944e-01
T4	-	1.0605911852062579e+02	1.0605911852062738e+02	1.06059118520625632823803211457878e+02	1.06059118520625632823803211457957e+02
T5	0	7.2164496600635175e-16	7.5495165674510645e-15	6.16297582203915472977912941627187e-32	7.90016463187644146923562152048337e-31

Table 4: Minimum and maximum final values found by the D and DD codes for the T functions,  $n = 10$ , and the optimal value.

Problem #	$f^*$	D		DD	
		min	max	min	max
T1	0	1.4601138701364973e-15	5.3033915333941507e-15	1.04077052403730597882313188466135e-31	4.78213141783335944624492452218446e-31
T2	0	8.1392709499333860e-09	3.2645285200950946e-07	7.38422049500142581209762558409305e-15	1.32699936552973002946751645600318e-12
T3	0	5.2624571367232420e-14	5.0598160442741924e-01	1.71947025434892416960837710713982e-30	5.05973651364335990653577129153149e-01
T4	-	5.6317033620503435e+02	5.8799776162067542e+02	5.63170336204998739130455870927457e+02	5.87997761620669625756282362037082e+02
T5	0	5.9396931817445875e-15	2.5079938126282286e-13	5.40801128383935827538118606277848e-31	1.10436674871052878285935887183704e-29

Table 5: Minimum and maximum final values found by the D and DD codes for the T functions,  $n = 50$ , and the optimal value.

Problem #	$f^*$	D		DD	
		min	max	min	max
T1	0	1.0620278236118440e-14	2.1619668381963696e-14	8.74714516122606297848555744313087e-31	1.96100321869372239665757142029997e-30
T2	0	1.3242105659618322e-08	3.4949545577428036e-07	7.36021573643631110286635269854721e-14	1.21516202980395497525178368262557e-12
T3	0	1.4343876061581251e-01	4.6214915356057096e-01	7.75403659008250178613729750667026e-02	4.62072399121541681882224700663063e-01
T4	-	2.2378539839989403e+03	2.3952676732459868e+03	2.23785398399887367684157263585940e+03	2.39526767324585797651230511115300e+03
T5	0	1.2307932450994485e-12	1.0875109146546436e-10	4.70805130485126127694652143932541e-29	7.21352438438372659396069204048122e-28

Table 6: *Minimum and maximum final values found by the D and DD codes for the T functions,  $n = 200$ , and the optimal value.*

For T1 and T5, the D and DD codes worked as we expected, reducing the final function values close to machine precision, around  $10^{-16}$  for D and  $10^{-32}$  for DD. On T2, the codes reduced the final function value to slightly larger values, but still close to the minimal value. However, the results for T3 and T4 were quite different, with the final function values varying quite widely over the 10 runs. For this reason, we decided to run both codes on T3 and T4 for  $n = 10$  and  $50$  with 1000 random starting points<sup>7</sup>, graphing the results with the DD values sorted and plotted with the corresponding D value for the same starting point, as shown in Figures 6.1 and 6.2 on the following page.

It seems that for both T3 and T4, there may be multiple local minima. For T3, we know the global minimum value is 0, and while the method sometimes finds this value in 10 runs for  $n = 10$  and  $n = 50$ , for  $n = 200$ , 10 runs are not enough. As can be seen in Figure 6.1, 1000 runs of T3 result in several plateaus, suggesting the presence of multiple local minima. For T4, we do not know the global minimum value, which changes with the number of variables. When  $n = 10$ , all 10 runs on T4 find the same value, and for  $n = 50$  and  $200$ , 10 runs are enough to sometimes find what may be the minimum value. However, over 1000 runs, it can be seen in Figure 6.2 on the following page that there are several plateaus, suggesting multiple local minima. Further investigation is required to prove there are in fact other local minima.

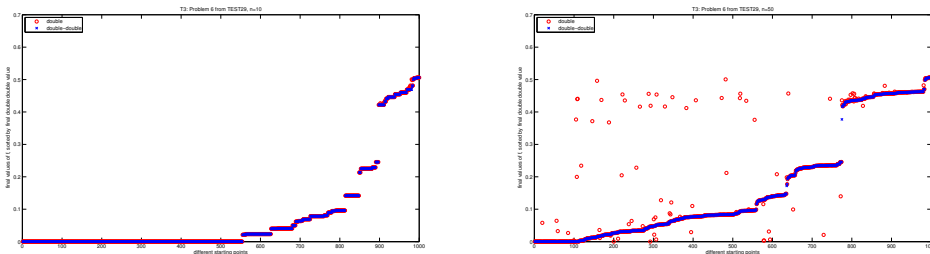


Figure 6.1: *T3 performance plots for  $n = 10$  (left) and  $n = 50$  (right), with the final DD function values sorted and plotted with the corresponding final D function value for 1000 random starting points. It appears that there are multiple local minima.*

<sup>7</sup>Tests on T3 and T4 for  $n = 200$  were not completed because of the extraordinarily long amount of time one run of the DD code takes to run before breakdown due to machine precision.

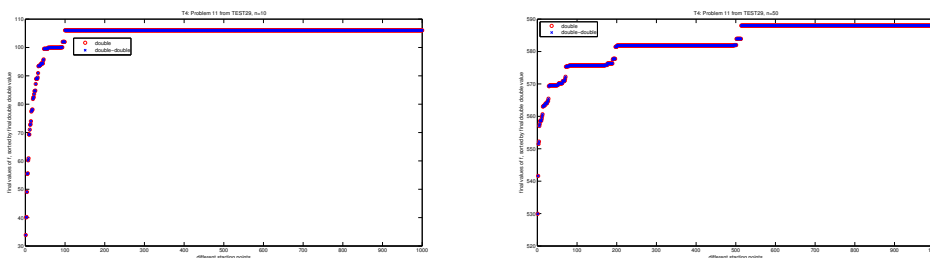


Figure 6.2:  $T_4$  performance plots for  $n = 10$  (left) and  $n = 50$  (right), with the final  $DD$  function values sorted and plotted with the corresponding final  $D$  function value for 1000 random starting points. It appears that there are multiple local minima.

### 6.3 Nesterov Functions Results

The three Nesterov functions are all very different from one another, and are thus affected differently by the implementation of higher precision.

The smooth function NCR-S is differentiable and we find, for  $n = 10$ , that using the D code we can reduce the gradient norm to about  $10^{-15}$ , while using the DD code, with only a few more iterations, we can reduce it to about  $10^{-161}$  before breakdown occurs. These gradient norm values correspond to final function values of approximately  $10^{-31}$  and  $10^{-323}$  for D and DD, respectively. As shown for  $n = 10$  in Table 7 on the next page, the D and DD codes generate function values that agree to double precision accuracy for many iterations, but eventually they start to diverge. By iteration 1988, they agree to only three digits, and by iteration 55025, they have no significant digits in common, although the order of magnitude is the same for both. The D code terminates at iteration 55036, while the DD code continues running until iteration 55102. The fact that the DD code makes so much improvement to the gradient norm with a relatively small increase in the number of iterations is evidence of superlinear convergence. The fourth and eighth columns show the gradient norm ratios  $\frac{\|\nabla \tilde{f}(x^{(j)})\|}{\|\nabla \tilde{f}(x^{(j-1)})\|}$ , and we see that indeed this ratio converges to zero, as required for superlinear convergence.

Figure 6.3 on page 21 and Figure 6.4 on page 22 shows the final function values found by both the D and DD codes for NCR-NS1 and NCR-NS2, respectively. The final function values found by the DD code have been sorted, and plotted with the corresponding function value the D code finds for the same starting point.

For the first nonsmooth variant NCR-NS1, the DD code overall performs better at converging to the global minimizer than the D code. For  $n = 2$ , the D and DD codes consistently reduce the final function value to about  $10^{-15}$  and  $10^{-30}$ , respectively, for all 1000 runs. For  $n = 3$ , the problem is much more difficult and BFGS already has trouble “tracking” the manifold  $M$ . As a result, for all starting points the D code is able to reduce the function value only to values in the range  $10^{-3}$  to  $10^{-8}$ ,



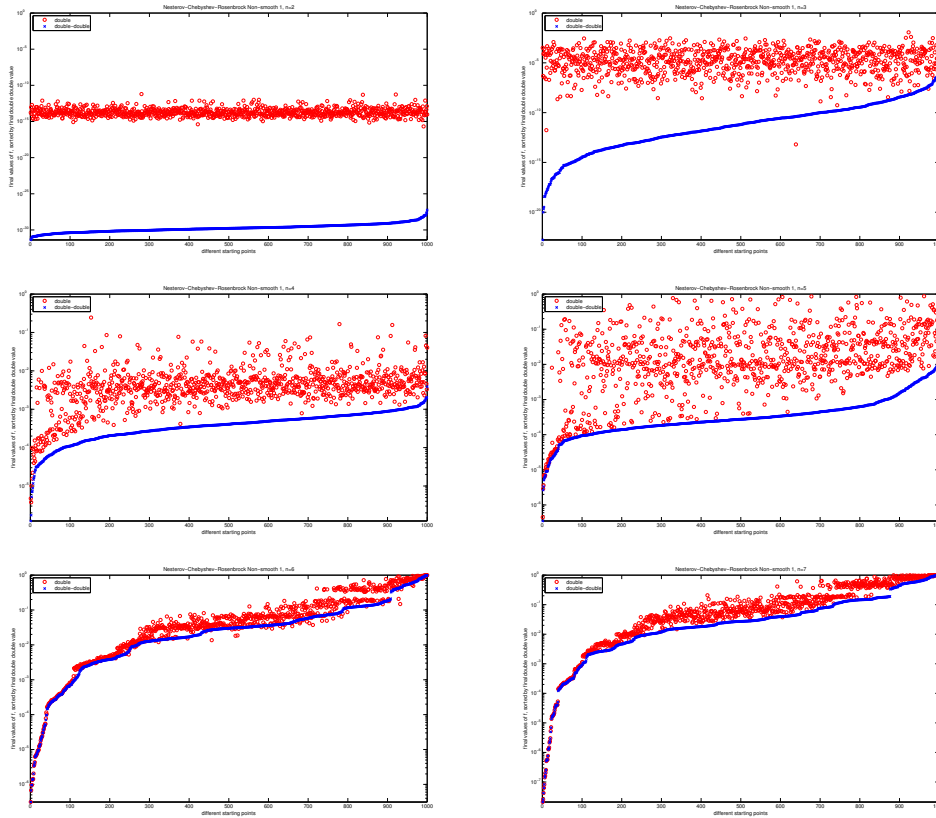


Figure 6.3: *NCR-NS1* performance plots, with the final *DD* function values sorted and plotted with the corresponding final *D* function value for 1000 starting points. Top left: Plot for  $n = 2$ . The *DD* code reduces the function value to as low as its precision allows. Top right: Plot for  $n = 3$ . The *DD* code overall does better than the *D* code, but with a much larger range. Middle left: Plot for  $n = 4$ . The *DD* code again does better than the *D* code, but final values are much closer. Middle right: Plot for  $n = 5$ . Both codes do equally well in some runs, and both codes have a large range of final values. Bottom left: Plot for  $n = 6$ . A limit is placed on the maximum number of iterations for  $n \geq 6$ . Neither code does very well, with little difference between the two. Bottom right: Plot for  $n = 7$ . Neither code does very well, with little difference between the two.

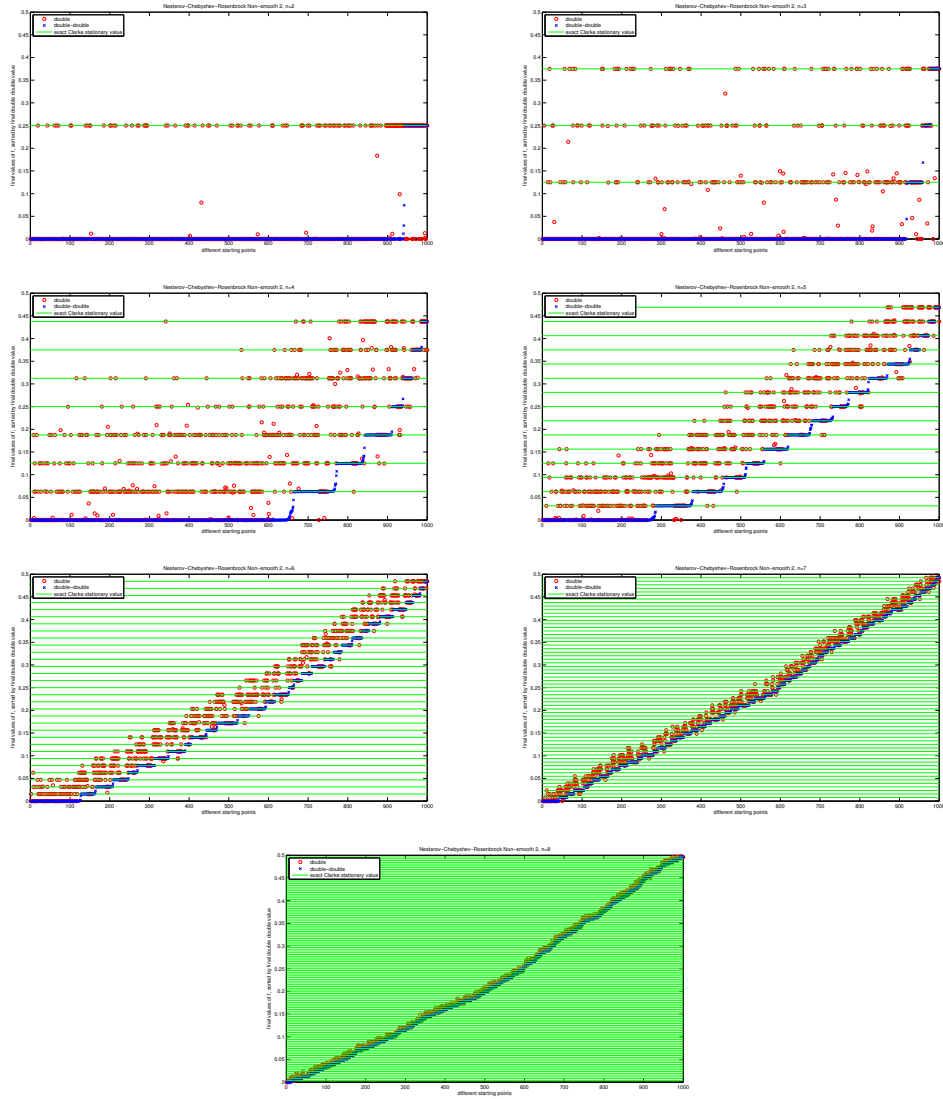


Figure 6.4: *NCR-NS2* performance plots, with the final *DD* function values sorted and plotted with the corresponding final *D* function value for 1000 starting points. Top left: Plot for  $n = 2$ . The *DD* code finds the minimizer much more often than the *D* code. Top right: Plot for  $n = 3$ . The *DD* code finds the minimizer most of the time, while the *D* code converges to the other Clarke stationary points more often. Top-Middle left: Plot for  $n = 4$ . The *DD* code finds the minimizer rather than the Clarke stationary points much more often than the *D* code. Top-Middle right: Plot for  $n = 5$ . The *DD* code finds a lower stationary point than the *D* code almost every time, but only finds the minimizer on a quarter of the runs. Bottom-Middle left: Plot for  $n = 6$ . Neither code finds the minimizer very often, but the *DD* code still finds a lower stationary point most of the time. Bottom-Middle right: Plot for  $n = 7$ . Neither code does very well in finding the minimizer a majority of the time, with minimal improvement of the *DD* code over the *D* code. Bottom: Plot for  $n = 8$ . Both codes perform practically the same.

but the DD code does much better, with final results ranging from  $10^{-6}$  to  $10^{-23}$ . When  $n = 4$ , both codes do much worse in converging to the global minimizer, but DD still performs better, with the D code final results ranging from  $10^{-1}$  to  $10^{-4}$ , while the DD code reduces the function value to values ranging from  $10^{-2}$  to  $10^{-6}$ . For  $n = 5$ , the D code final function values vary widely from 1 to  $10^{-6}$ , while the DD code performs slightly better, resulting in a tighter range of function values,  $10^{-2}$  to  $10^{-7}$ . In all the cases for  $n = 2, \dots, 5$ , the DD code does much better than the D code in almost all but a couple runs. The final function values form an “S”-shaped smooth curve, with a long, slowly increasing middle area, presumably reflecting the use of a uniform distribution for the starting points.

When  $n = 6$  and  $n = 7$ , the number of maximum iterations was limited to  $10^{n-1}$ . This was necessary in order to run the tests in a timely manner, as each run of the DD code would run for a couple of minutes before the machine precision causes termination of the code. Under these additional restrictions, the DD code has much less advantage over the D code. More runs of the D code produce the same if not better results. Also, the final function value curve loses its S-shape, looking more like a logarithmic curve. This indicates that although the D and DD codes are able to occasionally reduce the final function value to around  $10^{-7}$ , a majority of the time the algorithm is unable to reduce the final function value below  $10^{-3}$ . When the DD code was allowed to run with as many iterations as needed for  $n = 6$ , around three to four million iterations were required before the program terminated; however the final function value was still around  $10^{-7}$ . A complete run for NCR-NS1  $n = 7$  was not performed due to the limitations of the machine.

On the second nonsmooth variant NCR-NS2, the DD code does significantly better than the D code. As the plots show, starting from enough starting points BFGS “finds” all the Clarke stationary points, whose  $f$  values are marked on the plots by horizontal green lines. Intuitively, BFGS has trouble tracking the set  $S$  to the global minimizer because it cannot get past the corners, and as we can see, the algorithm ends up converging to these corners in both the D and DD implementations. For  $n = 2$  the DD code reduces the final function value to about  $10^{-32}$  most of the time, where the D version reduces it to about  $10^{-16}$  noticeably less often. In general, both codes do well in approaching the global minimizer under their respective precision. For  $n = 3$ , the DD code reduces the final function value to about  $10^{-31}$  a large majority of the time, however it also “finds” all the other Clarke stationary points. In contrast, the D code reduces the function value to  $10^{-15}$  only about two-thirds of the time, converging to the other Clarke stationary points much more often.

For  $n = 4$  we start to see a big difference in performance between the two codes. The DD code reduces the final function value to  $10^{-31}$  over half of the time, while the D code reduces it to  $10^{-15}$  only about a quarter of the time. The DD code still finds all the other Clarke stationary points, but the higher the corresponding function value, the less likely the DD code is to terminate there. The D code, on the other hand, converges to the other Clarke stationary points almost as often as the global minimizer. For  $n = 5$  the DD code reduces the final function value to approximately

$10^{-30}$  about a quarter of the time, but the D code only reaches  $10^{-14}$  a tenth of the time. The DD code finds the other Clarke stationary points decreasingly less often as the function value increases, while the D code finds most of the other Clarke stationary points almost as often as the global minimizer, and rarely converges to a lower final function value than the DD code.

Starting at  $n = 6$ , neither code does very well in finding the global minimizer. The DD code reduces the final function value to  $10^{-30}$  about a tenth of the time, while the D code reduces it to  $10^{-14}$  around 60 times in 1000 runs. In both codes, all other Clarke stationary points are found about as frequently as the global minimizer. However, on practically every run, the DD code reduces the final function value lower than the D code for the same starting point. For  $n = 7$ , although the DD code finds a lower final function value than the D code on almost every run, neither code does well in converging to the global minimizer a significant amount of times. By the time  $n = 8$ , there is almost no difference between the performance of the D and the DD code. Neither code successfully reduces the final function value to the global minimizer any more often than to any of the other Clarke stationary values.

On each test, the DD version is able to get past the corners to a lower point, possibly even the global minimizer, more often than the D code. However, occasionally the D code is able to reduce the final function value lower than the DD code, which we hypothesize happens because the D code gets lucky with rounding in the right direction, lower than the DD code does. However, this phenomenon occurs so seldom, it has almost no statistical significance.

To better demonstrate the ability of the DD code to reduce the final function value and get past the corners, where the D code could not, we looked at the function values for each iteration during some of the runs where a lower Clarke stationary point or the global minimizer was obtained by the DD code. These iteration plots for NCR-NS2 are shown in Figure 6.5 on the following page. In the graph where  $n = 3$ , although the DD code does not reach the global minimizer, it demonstrates its ability to get past the first corner and reach a lower Clarke stationary value. In the case where  $n = 5$ , the D code terminates after only 215 iterations at a Clarke stationary value, while the DD code continues on past several more corners, eventually reaching and terminating at the global minimizer. Finally, in the graph where  $n = 7$ , the D code again terminates very quickly after only 305 iterations, stopping at a Clarke stationary value. The DD code, however, continues on for several tens of thousands more iterations, eventually reaching a final function value of  $10^{-30}$ , where the code finally terminates.

In general, the DD runs for many more iterations than the D code before the program terminates. Consequentially, the DD code takes much longer to run. However, because of the huge effect different starting points have on the algorithm's ability and speed at finding the minimizer, it is difficult to find meaningful statistics on the time and number of iterations required. Some of the DD code runs took much longer than the corresponding D code, because of the longer time it takes to perform each



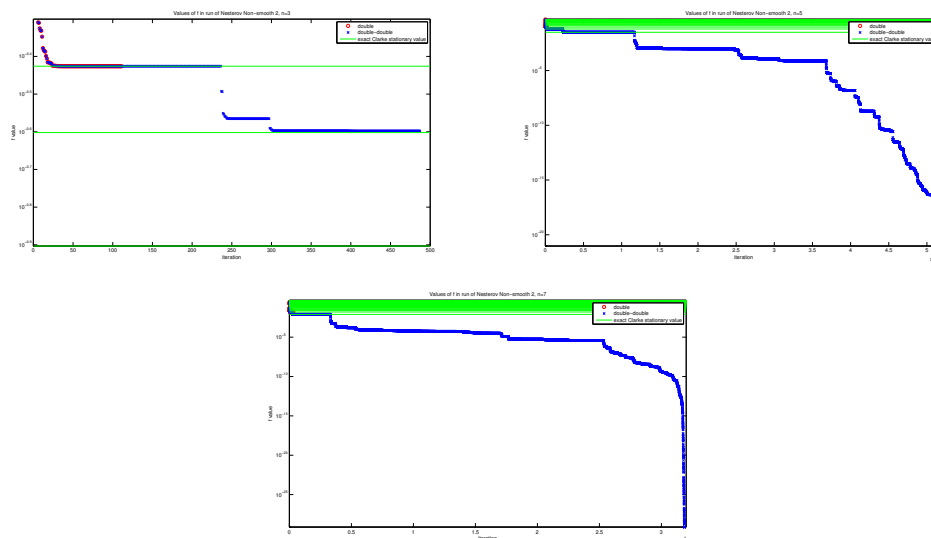


Figure 6.5: *Function values after each iteration of several runs of NCR-NS2.*

*Top left: Plot for  $n = 3$ . The DD code is able to find a smaller stationary point than the D code after several more iterations beyond the D code termination. Top right: Plot for  $n = 5$ . The DD code is able to get past several more corners than the D code. Bottom: Plot for  $n = 7$ . The DD code gets past many more corners than the D code.*

floating point operation in double-double precision, the larger number of iterations the DD code required to converge, and the longer time spent in the linesearch in each iteration. For example, in one run of NCR-NS1,  $n = 5$ , the D code took 7841 iterations, took an average of 2.54 expansion and/or bisection steps per line search, and took 0.04 seconds to run. The DD code took 342,830 iterations, took an average of 3.14 expansion and/or bisection steps per line search, and took 20.54 seconds to run. In addition, the D code has at most 25 expansion and/or bisection steps per line search, while the DD code took at most 72 expansion and/or bisection steps per line search. The DD code took approximately 44 times more iterations, but ran for about 514 times longer before terminating.

## 7 Conclusion

Using the BFGS code written in C by Skajaa [Ska10], we were able to implement double-double precision and test the hypothesis by Lewis and Overton [LO10] that BFGS is limited by machine precision, rather than the algorithm itself. In general, the DD code was able to reduce the final function values closer to the minimal value than the D code. In many cases, both codes were able to reduce the function essentially to the minimum value,  $10^{-16}$  for the D code and  $10^{-32}$  for the DD code.

---

Depending on the specific problem and the starting points chosen, the DD code was sometimes able to reduce the final function value significantly farther than the D code, particularly in the case of getting past Clarke stationary points in NCR-NS2.

As we saw on the Nesterov functions, the DD code performed much better overall than the D code for all three variations. For the smooth version, the DD code reduced the final function value much lower than the D code,  $10^{-323}$  compared to  $10^{-32}$ , without taking an exorbitant amount of time before terminating under machine precision. On the first nonsmooth variation, the DD code proved much better at reducing the final function value better than the D code for a small number of variables. For  $n = 6$  and larger, neither code performed well, and the DD version only did slightly better than the D code. However, the DD code took much longer to run, partly because of the extra iterations taken, and partly because each linesearch required more bisections and expansions before a sufficient stepsize was found.

For the second nonsmooth variation of Nesterov's function, the DD code proved much better at finding the minimizer and getting past the Clarke stationary points than the D code. When  $n = 2, \dots, 5$  the DD code found the minimizer more often than the other stationary points, and in general converged to a lower function value than the corresponding D code. For  $n = 6$  and larger, neither code does particularly well in finding the minimizer more than the other stationary points; however the DD code still does slightly better. By  $n = 8$ , however, there is almost no difference in the final function values of the two codes. With few exceptions, the double-double precision code proved much better at finding the minimizer and getting past the Clarke stationary points. However, for a larger number of variables and more difficult problems, much longer time is needed for the double-double precision code to run, and in some cases, an unrealistic amount of time is required before the code breaks down because of machine precision.

## Acknowledgements

I would like to thank my advisor, Professor Overton, for all the hard work he put in to helping me complete this thesis.

I would also like to thank Anders Skajaa for the critical role he played in translating the majority of the BFGS code into C, making my work far less than it would have been otherwise.

I also want to thank Andreas Kloeckner and Ben Stifter for the significant programming help they provided me on various occasions.

A big thanks also to Esteban Tabak for being my second reader on this thesis.

Finally, I would like to thank the following people for their small, but significant contributions:

Rachael Bailine  
Evan Herring  
Edwin Marte  
Smeet Merchant  
Ethan Price-Livingston  
Brian Sampson  
Serge Yegi

## References

- [Bai10] D.H. Bailey. A Fortran-90 based multiprecision system. *ACM Transactions on Mathematical Software*, 21(4):379–387, December 2010.
- [GO10] M. Gurbuzbalaban and M.L. Overton. On Nesterov’s nonsmooth Chebyshev-Rosenbrock functions. *Submitted to J. Nonlinear Analysis: Theory, Methods, Applications*, December 2010.
- [HMM04] M. Haarala, K. Miettinen, and M.M. Makela. New limited memory bundle method for large-scale nonsmooth optimization. *Optimization Methods and Software*, 19(6):673–692, 2004.
- [Kiw85] K.C. Kiwiel. *Methods of Descent for Nondifferentiable Optimization. Lecture Notes in Mathematics 1133*. Springer Verlag, 1985.
- [LDB<sup>+</sup>02] X. Li, J. Demmel, D. Bailey, G. Henry, Y. Hida, J. Iskandar, W. Kahan, A. Kapur, M.C. Martin, B. Thompson, and D.J. Yoo. Design, implementation, and testing of extended and mixed precision BLAS. *ACM Trans. Math. Software*, 2002.
- [Lew03] A.S. Lewis. Active sets, nonsmoothness and sensitivity. *SIAM Journal on Optimization*, 13:702–725, 2003.
- [LHB01] X. Li, Y. Hida, and D.H. Bailey. Algorithms for quad-double precision floating point arithmetic. In *15th IEEE Symposium on Computer Arithmetic*, pages 155–162, June 2001.
- [LO10] A.S. Lewis and M.L. Overton. *Nonsmooth Optimization via Quasi-Newton Methods*. August 2010.
- [LTS<sup>+</sup>02] L. Luksan, M. Tuma, M. Siska, J. Vlcek, and N. Ramesova. UFO 2002. Interactive system for universal functional optimization. Research report, Academy of Sciences of the Czech Republic, 2002.
- [NW06] J. Nocedal and S. Wright. *Numerical Optimization*. Springer, 2nd edition, 2006.
- [Ove01] M.L. Overton. *Numerical Computing with IEEE Arithmetic*. SIAM, 2001.
- [Pri91] D.M. Priest. Algorithms for arbitrary precision floating point arithmetic. In P. Kornerup and D. Matula, editors, *Proceedings of the 10th Symposium on Computer Arithmetic*, pages 132–143, Piscataway, NJ, 1991. IEEE Computer Society Press.
- [Ska10] A. Skajaa. Limited memory BFGS for nonsmooth optimization. Master’s thesis, Courant Institute of Mathematical Science, New York University, January 2010. <http://www2.imm.dtu.dk/~andsk/nsosite/webfiles/index.shtml>.