

Behavior of Limited Memory BFGS when Applied to Nonsmooth Functions and their Nesterov Smoothings

Azam Asl^{*}

Michael L. Overton[†]

June 15, 2020

Abstract

The motivation to study the behavior of limited-memory BFGS (L-BFGS) on nonsmooth optimization problems is based on two empirical observations: the widespread success of L-BFGS in solving large-scale smooth optimization problems, and the remarkable effectiveness of the full BFGS method in solving small to medium-sized nonsmooth optimization problems, based on using a gradient, not a subgradient, oracle paradigm. We first summarize our theoretical results on the behavior of the scaled L-BFGS method with one update applied to a simple convex nonsmooth function that is unbounded below, stating conditions under which the method converges to a non-optimal point regardless of the starting point. We then turn to empirically investigating whether the same phenomenon holds more generally, focusing on a difficult problem of Nesterov, as well as eigenvalue optimization problems arising in semidefinite programming applications. We find that when applied to a nonsmooth function directly, L-BFGS, especially its scaled variant, often breaks down with a poor approximation to an optimal solution, in sharp contrast to full BFGS. Unscaled L-BFGS is less prone to breakdown but conducts far more function evaluations per iteration than scaled L-BFGS does, and thus it is slow. Nonetheless, it is often the case that both variants obtain better results than the provably convergent, but slow, subgradient method. On the other hand, when applied to Nesterov's smooth approximation of a nonsmooth function, scaled L-BFGS is generally much more efficient than unscaled L-BFGS, often obtaining good results even when the problem is quite ill-conditioned. Summarizing, we find that although L-BFGS is often a reliable method for minimizing ill-conditioned smooth problems, when the condition number is so large that the function is effectively nonsmooth, L-BFGS frequently fails. This behavior is in sharp contrast to the behavior of full BFGS, which is consistently reliable for nonsmooth optimization problems. We arrive at the conclusion that, for large-scale nonsmooth optimization problems for which full BFGS and other methods for nonsmooth optimization are not practical, it is often better to apply L-BFGS to a smooth approximation of a nonsmooth problem than to apply it directly to the nonsmooth problem.

1 Introduction

We consider the unconstrained optimization problem

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

where the function f is convex but nonsmooth. By this we mean that it is not differentiable everywhere, and, typically, is not differentiable at minimizers.

^{*}Booth School of Business, University of Chicago

[†]Courant Institute of Mathematical Sciences, New York University

Classical approaches to optimization of convex nonsmooth functions generally require the method to have access to an oracle that, given $x \in \mathbb{R}^n$, returns the function value $f(x)$ and a subgradient $g \in \partial f(x)$. The oldest such method, the subgradient method of Shor, which dates to the 1970s, uses the iteration

$$x_{k+1} = x_k - t_k g_k, \text{ for some } g_k \in \partial f(x_k), \quad (1)$$

where $\{t_k\}$ is a pre-determined sequence of positive stepsizes. One well-known result states that, assuming f is convex and bounded below, and provided the steplengths $\{t_k\}$ are square-summable (that is, $\sum_{k=0}^{\infty} t_k^2 < \infty$, and hence the steps are “not too long”), but not summable (that is, $\sum_{k=0}^{\infty} t_k = \infty$, and hence the steps are “not too short”), then convergence of $f(x_k)$ to the minimal value of f must take place [NB01]. However, despite the strength of this theoretical result, it is well known that convergence to the optimal value is often very slow. This observation led to the development of other algorithms, particularly the bundle methods pioneered by Lemaréchal [Lem75] for nonsmooth convex functions in the mid-1970s, as well as the bundle methods of Kiwiel [Kiw85] for nonsmooth, nonconvex problems in the 1980s. As suggested by the name, at each iteration, bundle methods use subgradient information obtained at former iterates as well as at the current iterate x_k to generate the next iterate x_{k+1} . Convergence results are available for these methods, but the computational cost per iteration is significant for large n as computing x_{k+1} from x_k usually requires the solution of a quadratic program in n variables [BKM14, p. 306 and 313]. For more methods for nonsmooth optimization using the subgradient oracle, see [BKM14, BGK⁺20].

In this paper, we do not use the subgradient oracle paradigm. One reason for this is that, in practice, in the presence of rounding errors, it is often difficult, if not impossible, to determine whether the function f is differentiable at a given point x and hence to return a vector g which can be guaranteed to be a subgradient. A second is that since convex functions (and more generally, locally Lipschitz functions) are differentiable almost everywhere, there is no reason, at least in the absence of exact line searches, to suppose that a method will ever generate a point x_k where f is not differentiable. We therefore use the simpler paradigm that, given x , an oracle returns $f(x)$ and $g = \nabla f(x)$, if f is differentiable at x . Clearly, the subgradient and gradient oracles coincide when f is indeed differentiable at x . What if this is not the case? In theory, we need to assume that the oracle informs the method that the function is not differentiable and therefore a gradient cannot be provided. But in practice, the gradient oracle is simply implemented as if f is differentiable at x , breaking ties arbitrarily if necessary. For example, if $f(x)$ is defined as $\max(f_1(x), f_2(x))$, where f_1 and f_2 are smooth functions, the oracle returns $\nabla f_1(x)$ if $f_1(x) > f_2(x)$, $\nabla f_2(x)$ if $f_1(x) < f_2(x)$, and either one if $f_1(x) = f_2(x)$, a property that is difficult to determine anyway in the presence of rounding errors. A subgradient oracle might, in principle, return any convex combination of $\nabla f_1(x)$ and $\nabla f_2(x)$, but there is no reason for it to return anything other than $\nabla f_1(x)$ or $\nabla f_2(x)$, unless, for example, the intent is to return a “steepest descent” subgradient.

Using the gradient oracle instead of the subgradient oracle, we may ask the question: what can be said if we consider the ordinary gradient method,

$$x_{k+1} = x_k - t_k g_k, \text{ where } g_k = \nabla f(x_k), \quad (2)$$

where the steplength t_k is obtained, not from a predetermined sequence, but from a backtracking or Armijo-Wolfe line search? The answer is well known: it is often the case that x_k converges to a point \bar{x} where f is not differentiable and not minimal. See [AO20b] for a historical discussion and for a detailed analysis of an interesting special case. Note that in relevant illustrative examples, it is typical that f is differentiable at all iterates $\{x_k\}$,

and nondifferentiable only at the limit point \bar{x} . This demonstrates that the power of Shor’s subgradient method is not so much that it is defined even if f is not differentiable at x_k , as that the acceptable predetermined sequence of stepsizes $\{t_k\}$, unlike stepsizes generated by a line search, ensures convergence to a minimal value.

In the early 2000s, Burke, Lewis and Overton introduced the gradient sampling algorithm [BLO05] for nonsmooth, nonconvex optimization. This method uses the gradient, not the subgradient, oracle paradigm, and a convergence analysis for a rather general class of functions states that, with probability one, the method generates iterates $\{x_k\}$ where f is differentiable and that, if f is convex and bounded below, the function values $\{f(x_k)\}$ converge to the minimal value (more generally, that all cluster points of $\{x_k\}$ are Clarke stationary). However, the cost per iteration is prohibitive when the number of variables is large. See the survey paper [BCL⁺20] for more details, as well as recent work [CL20] introducing more efficient variants of the gradient sampling method.

As discussed by Lewis and Overton [LO13], the “full” BFGS quasi-Newton method is a very effective alternative choice for nonsmooth optimization, and its $O(n^2)$ cost per iteration is generally much less than the cost of the bundle or gradient sampling methods, but its convergence results for nonsmooth functions are limited to very special cases. It also uses the gradient, not the subgradient, oracle paradigm.

Since the limited memory variant of BFGS [LN89] (L-BFGS) costs only $O(n)$ operations per iteration, like the gradient and subgradient methods, it is natural to ask whether L-BFGS could be an effective method for nonsmooth optimization. This is the topic of this paper, which is organized as follows. In §2, we summarize our theoretical results. Then in §3, we give extensive experimental results. Although the methods we discuss are applicable to nonconvex problems, we restrict our discussion to the convex case for simplicity, focusing on a difficult nonsmooth problem of Nesterov as well as eigenvalue optimization problems, including those arising from semidefinite programming formulations of the max cut and matrix completion problems. We also consider Nesterov’s smooth approximations to these nonsmooth problems. We make some concluding remarks in §4.

2 Limited Memory BFGS for Nonsmooth Optimization in Theory

We begin this section by defining the concept of an Armijo-Wolfe line search. Then we discuss the full BFGS method on which L-BFGS is based, along with its properties, before discussing our results for the L-BFGS method.

2.1 Armijo-Wolfe Line Search

The BFGS and L-BFGS methods rely on an Armijo-Wolfe line search (also often known as a “weak Wolfe” line search). Let c_1 and c_2 be parameters (the Armijo parameter and the Wolfe parameter respectively) satisfying $0 < c_1 < c_2 < 1$. Let $x_k \in \mathbb{R}^n$ be a given iterate where f is differentiable and let $d_k \in \mathbb{R}^n$ be a descent direction for f from x_k , that is, with $\nabla f(x_k)^T d_k < 0$. We say that a positive steplength t_k satisfies the Armijo condition if

$$f(x_k + t_k d_k) \leq f(x_k) + c_1 t_k \nabla f(x_k)^T d_k, \quad (3)$$

and that t_k satisfies the Wolfe condition if f is differentiable at $x_k + t_k d_k$ and

$$\nabla f(x_k + t_k d_k)^T d_k \geq c_2 \nabla f(x_k)^T d_k. \quad (4)$$

The Armijo condition imposes a “sufficient decrease” in the value of $f(x_k + t_k d_k)$ compared to $f(x_k)$, while the Wolfe condition imposes a “sufficient increase” in the directional derivative

$\nabla f(x_k + t_k d_k)^T d_k$ compared to the negative value $\nabla f(x_k)^T d_k$. If f is bounded below on the ray $\{x_k + t d_k : t > 0\}$, then it is known that points t_k satisfying these conditions exist under various assumptions on f , such as f being convex (but not necessarily smooth). For further discussion, including specification of a bracketing line search algorithm based on bisection and doubling to compute t_k , see [LO13].

While it is often recommended in the context of general smooth nonlinear optimization to set the Armijo parameter to a rather small value such as 10^{-4} [NW06, p. 62], we note that for satisfactory complexity results for the gradient method applied to smooth, strongly convex functions, the Armijo parameter must *not* be too small [BV04, p.466–468]. An analysis of the gradient method using an Armijo-Wolfe line search applied to a nonsmooth function that we discuss later in this paper (see (9)) was given in [AO20b]. It was shown that, in this case, the success or failure of the gradient method depends critically on the Armijo parameter being sufficiently small, and experiments applying the gradient method to another nonsmooth function (10) confirm the importance of this issue [AO20b, Fig. 6]. On the other hand, as we will see in §2.3, the choice of Armijo parameter is less critical to the success or failure of L-BFGS when applied to the same function (9).

2.2 Full BFGS

The Broyden-Fletcher-Goldfarb-Shanno (BFGS) method, proposed independently by these four authors in 1970, is a quasi-Newton method that maintains an approximation H_k to the inverse of the Hessian $\nabla^2 f(x_k)$. Let an initial iterate $x_0 \in \mathbb{R}^n$ and an initial positive definite matrix H_0 be given. The BFGS method is defined by the following iteration. For $k = 0, 1, 2, \dots$, let

$$d_k = -H_k \nabla f(x_k)$$

$$x_{k+1} = x_k + t_k d_k \text{ where } t_k > 0 \text{ satisfies (3), (4)}$$

$$s_k = t_k d_k \tag{5}$$

$$y_k = \nabla f(x_{k+1}) - \nabla f(x_k) \tag{6}$$

$$\rho_k = \frac{1}{s_k^T y_k}$$

$$H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) + \rho_k s_k s_k^T \tag{7}$$

Equation (7) is called the BFGS update. It is easy to check that the Wolfe condition ensures that $s_k^T y_k > 0$, and it is well known that the positive definiteness of H_{k+1} follows as a consequence. It is also clear that H_{k+1} can be computed in $O(n^2)$ operations.

The most important convergence property of BFGS is due to Powell [Pow76]. It states that if f is twice continuously differentiable and strongly convex on the level set $S = \{x : f(x) \leq f(x_0)\}$, then the sequence $\{x_k\}$ converges to the global minimizer. Furthermore, the convergence theory of Dennis and Moré states that if we further assume that the Hessian of f , $\nabla^2 f(x)$, is Lipschitz continuous at the global minimizer, then the rate of convergence is superlinear. See [NW06] for more details.

For many years, BFGS was not considered as a possible stand-alone method for non-smooth optimization, although Lemaréchal briefly mentioned this possibility in a 1982 technical report. Instead, Lemaréchal and others focused on using the BFGS update to provide second-order information to bundle methods [BKM14, p. 313]. In contrast, Lewis and Overton [LO13] advocated using the pure BFGS method, with no bundle method component, as a practical method for nonsmooth optimization, using the gradient, not the subgradient, oracle paradigm. The Wolfe condition (4) ensures in theory that f is differentiable at all iterates, although in practice, as already noted, it is difficult to check whether or not f is

differentiable at a given point. If f is not differentiable at x , then the gradient oracle may return very different gradients at points that are close to x . The consequence is that, in the nonsmooth case, BFGS builds a very ill-conditioned inverse “Hessian” approximation, with some tiny eigenvalues corresponding to “infinitely large” curvature in the directions defined by the associated eigenvectors. This phenomenon, illustrated in [LO13, Fig. 4], is apparently what makes BFGS so effective for nonsmooth optimization. Remarkably, the condition number of the inverse Hessian approximation often reaches 10^{16} before the method breaks down numerically, usually because, as a consequence of rounding error, the line search is unable to return t_k satisfying the Armijo condition even though, in principle, an Armijo-Wolfe step exists. Convergence of $f(x_k)$ to the minimal value of f often appears to be linear (not superlinear, as in the smooth case).

Although empirically BFGS works very well when f is nonsmooth, convergence results are limited to a few special cases. The following results hold for all x_0 as long as f is differentiable at x_0 and for all positive definite H_0 . We use $x^{(i)}$ to denote the i th entry of the vector x .

- $n = 1$ with $f(x) = |x|$: the sequence generated by BFGS converges linearly to zero and is related to a certain binary expansion of the starting point [LO13].
- $f(x) = |x^{(1)}| + \sum_{i=2}^n x^{(i)}$: eventually a direction is identified on which f is unbounded below [XW17]; see also [LZ15].
- $f(x) = \|x\|_2$: iterates converge to the origin [GL18]. This is a special case of a more general result whose proof is based on Powell’s theorem mentioned above.

As far as we know, even the case $f(x) = |x^{(1)}| + (x^{(2)})^2$ remains open!

BFGS has been used successfully in many practical applications of nonsmooth optimization including the design of fixed-order controllers for linear dynamical systems with input and output, shape optimization for spectral functions of Dirichlet-Laplacian operators and condition metric optimization. For more details, see [LO13, p. 159–160]. Software is available in the HANSO package.¹

Finally, BFGS has also proved useful for optimization problems with nonsmooth constraints. Consider problems of the form

$$\begin{aligned} \min \quad & f(x) \\ \text{subject to } & c_i(x) \leq 0, \quad i = 1, \dots, p \end{aligned}$$

where f and c_1, \dots, c_p may not be differentiable at local minimizers. A successive quadratic programming (SQP) method based on BFGS was introduced by [CMO17] to solve problems of this form, and applied to problems in static output feedback control design that involve spectral radius and pseudospectral radius constraints. Although there are no theoretical results, it is typically much more efficient in practice than an SQP gradient sampling method [CO12] which does have convergence results. Software is available in the GRANSO package.²

2.3 Limited Memory BFGS

Unlike the full BFGS method, L-BFGS does not store a matrix approximating the inverse of the Hessian $\nabla^2 f(x)$. Instead, it uses an implicit approximation. Let m be a small integer. Instead of using (7) to update a stored matrix H_k , the matrix H_{k+1} is implicitly defined by applying m BFGS updates successively, starting with an initial matrix $H_{k+1}^{(0)}$ and using

¹<http://www.cs.nyu.edu/~overton/software/hanso/>

²<http://www.timmitchell.com/software/granso/>

the updates defined by the pairs $(s_j, y_j), j = k - m + 1, \dots, k$, defined in (5) and (6). This method is called L-BFGS- m . We consider two variants. In the *scaled* variant, with

$$H_{k+1}^{(0)} = \frac{s_k^T y_k}{y_k^T y_k} I, \quad (8)$$

a scaling of the identity matrix often known as Barzilai-Borwein scaling [BB88] is used to initialize the implicit updating at iteration k . In the *unscaled* variant, $H_{k+1}^{(0)}$ is set to the identity matrix. We also refer to these two different variants of L-BFGS- m as *with* and *without scaling*, respectively. Substantial numerical experience with applying L-BFGS- m to minimize smooth functions shows that the use of scaling is strongly preferred. For more details on L-BFGS, including efficient implementation, see [LN89, NW06].

Our theoretical analysis of the behavior of limited memory BFGS uses the scaled variant with $m = 1$: H_{k+1} is defined by applying just one BFGS update to (8). This method, L-BFGS-1, is sometimes called *memoryless BFGS* [NW06, p. 180]. We focus on the nonsmooth function

$$f(x) = a|x^{(1)}| + \sum_{i=2}^n x^{(i)} \quad (9)$$

with $a \geq \sqrt{n-1}$. This function is obviously unbounded below, but it is bounded below along any steepest descent direction $d = -\nabla f(x)$. One advantage of studying this function is its simplicity, but another is that it is easy to determine whether a method succeeds or fails when it is applied to (9): a method is successful only if it generates a sequence of function values $f(x^{(k)}) \rightarrow -\infty$ or identifies a direction d_k on which $\{f(x_k + td_k) : t > 0\}$ is unbounded below. Otherwise, since Armijo-Wolfe steps always exist along directions on which f is bounded below, the sequence of function values $f(x^{(k)})$ is well defined, bounded below, and must converge to a non-optimal finite value. The following theorem shows that this last scenario occurs when the Armijo parameter is not sufficiently small compared to the parameter a defining the function definition in (9). The proof is given in [AO20a, §3.2].

Theorem 1 *Suppose f is defined by (9) with $a \geq 2\sqrt{n-1}$. Set x_0 to any point with $x_0^{(1)} \neq 0$. If the Armijo parameter c_1 is chosen so that*

$$\frac{1 - c_1}{c_1}(n - 1) < a^2 + a\sqrt{a^2 - 3(n - 1)}$$

holds, then the scaled L-BFGS-1 method is well defined, in the sense that Armijo-Wolfe steps always exist, but fails in the sense that $f(x_k)$ is bounded below as $k \rightarrow \infty$.

Note that $a \geq 2\sqrt{n-1}$ implies

$$a^2 + a\sqrt{a^2 - 3(n - 1)} \geq 4(n - 1) + 2\sqrt{n - 1}\sqrt{n - 1} = 6(n - 1).$$

So, if $c_1 > 1/7$, the failure condition holds. It was proved in [AO20b] that, when applied to (9), the gradient method with an Armijo-Wolfe line search fails in the same sense if the stronger condition

$$\frac{1 - c_1}{c_1}(n - 1) < a^2$$

on the Armijo parameter holds. So, surprisingly, in some cases scaled L-BFGS-1 fails although the gradient method succeeds in generating $f(x_k) \rightarrow -\infty$.

If the specific Armijo-Wolfe bracketing line search given in [LO13] is used, we obtain a failure condition for scaled L-BFGS-1 that is *independent* of the Armijo parameter. The proof of the next result is given in [AO20a, §3.3].

Theorem 2 Suppose f is defined by (9) with $a \geq 2\sqrt{n-1}$. Set x_0 to any point with $x_0^{(1)} \neq 0$. If scaled L-BFGS-1 is implemented using the Armijo-Wolfe bracketing line search of [LO13], the method is well defined in the sense that the line search always returns an Armijo-Wolfe steplength t_k , but fails in the sense that $f(x_k)$ is bounded below as $k \rightarrow \infty$.

In fact, in the experiments reported in [AO20a] we observe that $a \geq \sqrt{3(n-1)}$ suffices for the method to fail. Furthermore, experiments indicate that, for any number of updates m , there is a threshold for the parameter a beyond which scaled L-BFGS- m always fails in the same sense: $f(x_k)$ is bounded below as $k \rightarrow \infty$. This threshold value for a increases with m .

3 Limited Memory BFGS for Nonsmooth Optimization in Practice

In this section we carry out extensive experiments applying L-BFGS- m , with and without scaling, using the Armijo-Wolfe bracketing line search of [LO13], to several challenging nonsmooth examples. All of the functions we consider from now on are bounded below. In most cases, in assessing whether a method succeeds or fails, we compare the final computed result with the known optimal value. We also compare the results with those obtained by the full BFGS method using the same line search. We used the HANSO package which implements both full and limited-memory BFGS, with the default settings for the Armijo and Wolfe parameters: $c_1 = 10^{-4}$ and $c_2 = 0.5$. In many cases, we also compared the results to those obtained by the subgradient method with predetermined stepsizes $t_k = 1/k$. All methods, including the subgradient method, are implemented using the gradient oracle paradigm discussed in §1: no attempt is made to determine whether the objective function is differentiable at a given iterate x_k . Instead of using a stopping criterion such as discussed in [LO13, p. 159], we run each method until it reaches a maximum number of function evaluations or iterations, or the line search fails to find an acceptable step (either because a limit on the number of bisections in the line search is exceeded or because of rounding errors). If the method terminates because of failure in the line search when $f(x_k)$ is not a good approximation to the optimal value, we say that the method “breaks down”. We find, in general, that this behavior is typical for L-BFGS, but not for full BFGS. For this reason we also compare these methods on *smoothed* versions of the objective functions.

We begin with a difficult nonsmooth problem devised by Nesterov, and then we go on to consider eigenvalue optimization and semidefinite programming problems.

3.1 Nesterov’s Les Houches Problem

We now consider a nonsmooth function introduced by Nesterov. The function is defined by

$$f(x) = \max\{|x^{(1)}|, |x^{(i)} - 2x^{(i-1)}|, i = 2, \dots, n\}. \quad (10)$$

Let $\hat{x}^{(1)} = 1, \hat{x}^{(i)} = 2\hat{x}^{(i-1)} + 1, i = 2, \dots, n$. Then $f(\hat{x}) = 1 = f(\mathbf{1})$, where $\mathbf{1} = [1, \dots, 1]^T$, although $\|\hat{x}\|_\infty \approx 2^n$ and $\|\mathbf{1}\|_\infty = 1$, so the level sets of f are very distorted. The minimizer is $\mathbf{0} = [0, \dots, 0]^T$ with $f(\mathbf{0}) = 0$. We call this Nesterov’s “Les Houches” function [Nes16].

In Figure 1, we show results for L-BFGS-1 (memoryless BFGS), with and without scaling, as well as full BFGS and the subgradient method, for the Les Houches problem with $n = 500$. We display all the function values that were computed, including those computed in the bracketing line search. We put a limit of 10,000 function evaluations on each method. The starting point x_0 (used by all the methods) was drawn randomly from the ball of radius 0.1 centered at the vector of all ones, using the standard normal distribution.

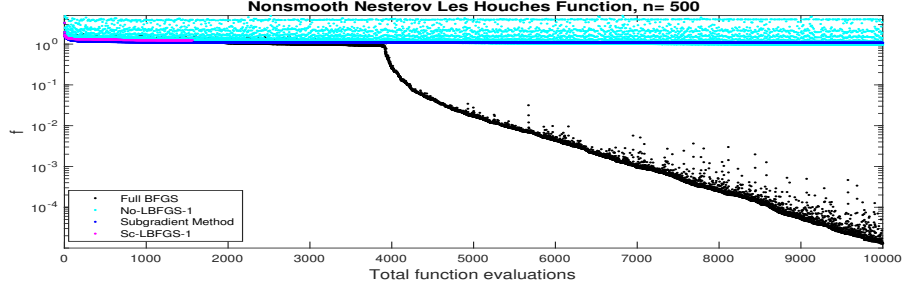


Figure 1: Comparing full BFGS, L-BFGS-1 with and without scaling and the subgradient method on the nonsmooth Les Houches problem (10) with $n = 500$. Here and below, “No-LBFGS” and “Sc-LBFGS” refer to the methods without scaling and with scaling respectively.

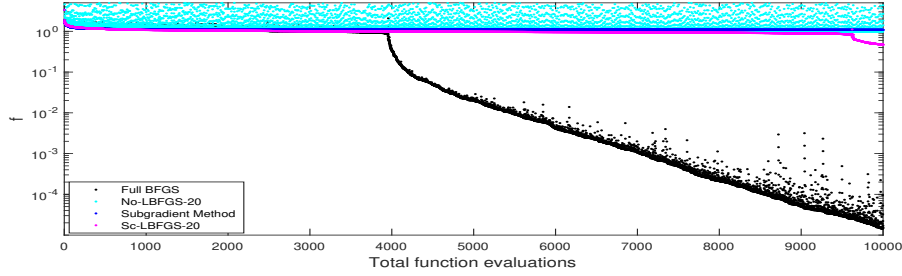


Figure 2: Comparing full BFGS, L-BFGS-20 with and without scaling and the subgradient method on the nonsmooth Les Houches problem (10) with $n = 500$.

We see from Figure 1 that scaled L-BFGS-1 (magenta dots) breaks down, with failure in the line search, after fewer than 2000 function evaluations. In contrast, unscaled L-BFGS-1 (cyan) runs for the full 10,000 function evaluations. However, its scattered plot indicates that the method performs many function evaluations per iteration in the line search, indicating that, not surprisingly given its name, the search directions it generates are not well scaled. Despite this, the method obtains a somewhat lower answer than the subgradient method (dark blue). Full BFGS (black) performs much better than any of the other methods, reducing the function value to about 10^{-5} (recall that the optimal value is zero). It is interesting to note that its convergence rate picks up rapidly right after it has lowered the function value down to 1. We do not know the reason for this.

We now increase the number of updates m from 1 to 20 and repeat this experiment: see Figure 2. Unlike scaled L-BFGS-1, scaled L-BFGS-20 does not quit early, and furthermore it also demonstrates a suddenly faster convergence rate toward the end of the experiment similar to that of full BFGS (although the final answer it obtains is not nearly as accurate as full BFGS). It obtains a function value of size 0.47, whereas unscaled L-BFGS-20 gets a final answer of about 0.998 and the subgradient method obtains 1.083. In this experiment, increasing the number of updates from 1 to 20 enhanced the performance of scaled L-BFGS far more than it did for unscaled L-BFGS.

The conclusion from these experiments is that with a small m , unscaled L-BFGS- m performs better and with a larger m it is the scaled variant which performs better. However, neither method performs nearly as well as full BFGS.

3.2 Smoothed Versions of Nesterov’s Les Houches Problem

Since L-BFGS- m performed poorly on (10), we consider instead applying it to a smoothed version. Let

$$A = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ -2 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & -2 & 1 \end{bmatrix}.$$

Then, (10) is equivalent to

$$f(x) = g(Ax), \quad (11)$$

where $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is defined by

$$g(y) = \max\{|y^{(i)}| : i = 1, 2, \dots, n\}. \quad (12)$$

Consider the Nesterov smoothing [Nes05, Van19] of the vector-max problem (12):

$$g_\mu(y) = \mu \log \sum_{i=1}^n (e^{y^{(i)}/\mu} + e^{-y^{(i)}/\mu}) - \mu \log(2n). \quad (13)$$

Without the constant term $-\mu \log(2n)$, this function is sometimes known as the softmax function [MRT12, p. 205]. The unique minimizer of $g_\mu(y)$ is $y_\mu^* = \mathbf{0}$ with $g_\mu^* = 0$. Since A is full-rank (although one of its singular values converges to zero as $n \rightarrow \infty$), via (11) we know that the unique minimizer of $f_\mu(x)$, x_μ^* , is also $\mathbf{0}$, with the same optimal value $f_\mu^* = 0$, and it can be verified that

$$\|\nabla^2 f_\mu(\mathbf{0})\|_2 = \frac{1}{n\mu} \|A\|_2^2. \quad (14)$$

We follow the standard approach [BV04, Sec 9.1] to defining the condition number of the strongly convex function f_μ as

$$\kappa(f_\mu) = \left(\max_{x \in S} \|\nabla^2 f_\mu(x)\|_2 \right) \left(\max_{x \in S} \|(\nabla^2 f_\mu(x))^{-1}\|_2 \right) \quad (15)$$

where $S = \{x : f(x) \leq f(x_0)\}$. For small μ , the second factor is enormous as all eigenvalues of $\nabla^2 f_\mu(x_0)$ are tiny. Using (14) as a lower bound for the first factor we conclude that, for small μ ,

$$\kappa(f_\mu) \gg \frac{1}{\mu}.$$

We now report on experiments we conducted applying full BFGS and L-BFGS, with and without scaling, to the smooth function f_μ , with $n = 500$ as before. All of the methods start from the same initial point used earlier. The left panel of Figure 3 shows the final function value computed by full BFGS and L-BFGS-1 with and without scaling as a function of the smoothing parameter μ using a log-log scale. Let us focus first on the results for full BFGS (black circles). BFGS always finds a solution with magnitude smaller than 10^{-15} , even for a very small μ , when the function is extremely ill conditioned. This is a remarkable property of full BFGS: its accuracy does not deteriorate significantly³ as the condition number $\kappa(f_\mu)$ of the smoothed problem blows up with $\mu \rightarrow 0$. In fact, when μ is sufficiently small, say $\mu = 10^{-16}$ (approximately the rounding unit in IEEE double precision used by MATLAB), the smoothed problem is essentially equivalent to the original nonsmooth problem when

³Surprisingly, the accuracy increases somewhat as μ decreases, but this is at the level of rounding errors and could perhaps be explained by a rounding error analysis. Certainly the scatter at the bottom left corner of the left panel of Figure 3 is a consequence of rounding error.

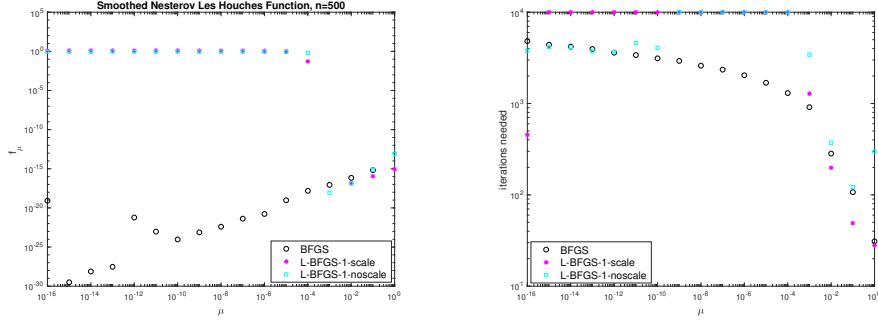


Figure 3: Comparing BFGS and L-BFGS-1 with and without scaling on the smoothed Les Houches function (13) for $n = 500$. The left panel shows the final function value and the right panel shows the iteration count, both as a function of the smoothing parameter μ . The maximum number of iterations is set to 10^4 .

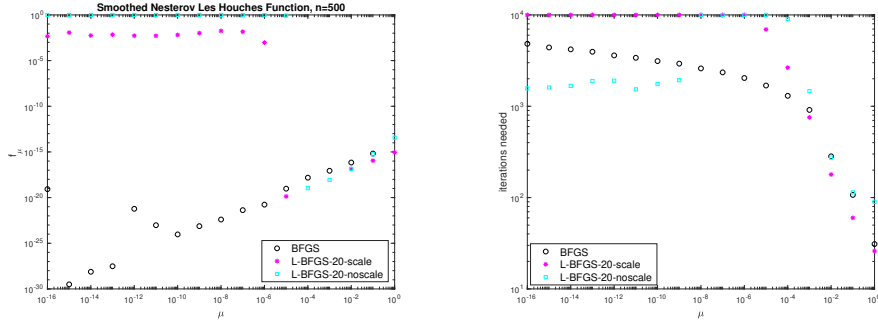


Figure 4: Comparing BFGS and L-BFGS-20 with and without scaling on the smoothed Les Houches function (13) for $n = 500$. The left panel shows the final function value and the right panel shows the iteration count, both as a function of the smoothing parameter μ . The maximum number of iterations is set to 10^4 .

rounding errors are taken into account, so the left panel shows the transition of the accuracy of full BFGS from smoothed variants of the problem to the limiting nonsmooth problem. The right panel shows the number of iterations that were required, again as a function of the smoothing parameter μ and again using a log-log scale. The maximum number of iterations (not function evaluations) was set to 10^4 for each μ . Remarkably, we see that the number of iterations required for full BFGS to accurately minimize f_μ does not significantly increase as $\mu \rightarrow 0$, even though the condition number $\kappa(f_\mu)$ blows up as μ decreases to zero, and the number required for the effectively nonsmooth instance $\mu = 10^{-16}$ is not much more than the number required for much better conditioned smoothed problems arising from moderate values of μ .

The results for L-BFGS-1 are very different. Unscaled L-BFGS-1 (cyan squares) finds an accurate answer for $\mu \geq 10^{-3}$, but the number of iterations required increases rapidly as μ is decreased further so the iteration limit is reached for μ ranging from 10^{-4} to 10^{-9} . However, starting with $\mu = 10^{-10}$, unscaled L-BFGS-1 breaks down before reaching the maximum number of iterations. The behavior of scaled L-BFGS-1 (magenta asterisks) is similar except that it breaks down only for $\mu = 10^{-16}$. Note that, since we are displaying the number of iterations, not the number of function evaluations, the performance of unscaled L-BFGS-1 looks better than it really is: the scaled version is computing substantially fewer function evaluations per line search.

When we increase the number of updates to $m = 20$, scaled L-BFGS reacts much better than unscaled L-BFGS; see Figure 4. Unscaled L-BFGS-20 finds accurate answers for $\mu \geq 10^{-4}$ before hitting the iteration limit, while the scaled version does so for $\mu \geq 10^{-5}$. Furthermore, when the maximum iteration limit is reached, scaled L-BFGS-20 achieves an answer of magnitude $\approx 10^{-2}$, whereas unscaled L-BFGS-20 is still giving an answer of magnitude $\approx 10^0$, similar to unscaled L-BFGS-1. However, overall, both are still doing poorly compared to full BFGS.

Our conclusions from this subsection are consistent with the generally accepted wisdom concerning L-BFGS. For smooth problems, even very ill-conditioned ones, it is best to use the scaled version of L-BFGS, and choosing the number of updates m to be larger rather than smaller gives better performance, although, in contrast to full BFGS, the number of iterations required increases significantly with the conditioning of the problem. Again in contrast with full BFGS, when the ill conditioning increases to the nonsmooth limit implicit in consideration of rounding errors, scaled L-BFGS generally fails to converge to an optimal solution. However, for the smooth but ill-conditioned problems considered in this subsection, unscaled L-BFGS offers no advantage compared to scaled L-BFGS. The most important conclusion is that, while applying full BFGS directly to nonsmooth problems works remarkably well, this is not the case for L-BFGS; at least for the Les Houches problem, it is far preferable to apply scaled L-BFGS to a smooth approximation of the nonsmooth problem.

3.3 Max Eigenvalue Problem

Let S^N denote the space of $N \times N$ real symmetric matrices, and let $\mathcal{A} : S^N \rightarrow \mathbb{R}^n$ denote a linear operator acting on X as follows:

$$\mathcal{A}X = \begin{bmatrix} \langle A_1, X \rangle \\ \vdots \\ \langle A_n, X \rangle \end{bmatrix}, \quad (16)$$

with $A_i \in S^N$ for $i = 1, \dots, n$. Its adjoint operator; $\mathcal{A}^T : \mathbb{R}^n \rightarrow S^N$, is defined by

$$\mathcal{A}^T y = \sum_{i=1}^n y^{(i)} A_i. \quad (17)$$

The Max Eigenvalue problem is to minimize the function

$$f(y) = \lambda_{\max}(C - \mathcal{A}^T y), \quad (18)$$

where $C \in S^N$ and $\lambda_{\max} : S^N \rightarrow \mathbb{R}$ denotes largest eigenvalue of its argument. It is well known that λ_{\max} is a convex function on S^N . Early papers on eigenvalue optimization include [Ove88].

Assuming the maximum eigenvalue of $C - \mathcal{A}^T y$ is simple, the gradient of f is

$$\nabla f(y) = -\mathcal{A}(qq^T) = -[q^T A_1 q, \dots, q^T A_n q]^T.$$

Following the gradient oracle paradigm discussed in §1, we make no attempt to estimate whether or not the maximum eigenvalue is simple at a given iterate. However, at optimal solutions, we generally expect that $C - \mathcal{A}^T y$ has a multiple largest eigenvalue and hence f is not differentiable. As is well known, eigenvalue optimization problems are instances of semidefinite programs, and hence small problems can be solved using CVX [GB14].

Using the standard normal distribution, we generated a random instance of this problem, defining \mathcal{A} and C with $N = 50$ and $n = 49$. Figure 5 shows the performance of full BFGS,

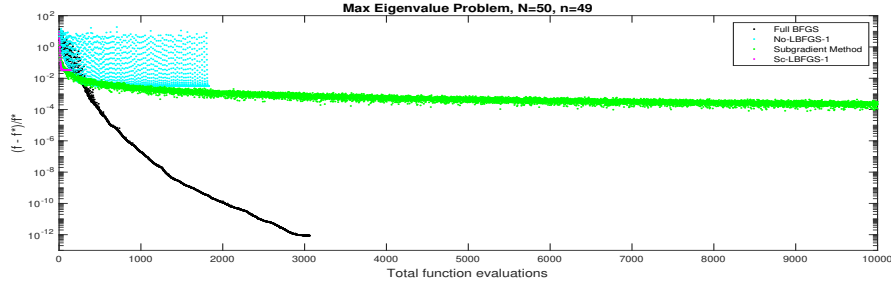


Figure 5: Comparing BFGS, L-BFGS-1 with and without scaling and the subgradient method on a randomly generated Max Eigenvalue problem (18) with $N = 50$ and $n = 49$.

SDPT3	BFGS	Sc L-BFGS-1	No L-BFGS-1	subgradient
7.82702970305352	7.82702970306035	8.08455876518360	7.85155000878711	7.82953885641783
7.82702970305349	7.82702970306035	8.08455876518359	7.85155000044960	7.82746561454846
7.82702970305348	7.82702970306034	8.08197715145863	7.85154996050940	7.82673229360627
7.82702970305346	7.82702970306031	8.05541534062475	7.85141043900471	7.82472408900949
7.82702970305334	7.82702970306017	7.84362627205676	7.69075549655739	7.82286893229481
7.70350538019538	7.70350432059448	7.56258926523925	7.48734455288558	7.70188538367848

Table 1: Top 6 eigenvalues of $C - \mathcal{A}^T y$ for Max Eigenvalue problem (18) where y is computed by SDPT3, full BFGS, scaled/unscaled L-BFGS-1 and the subgradient method for a randomly generated problem with $N = 50$ and $n = 49$. The optimal multiplicity is 5.

L-BFGS-1 with and without scaling, and the subgradient method (as before with $t_k = 1/k$) for minimizing (18). As earlier, we display all the function values that were computed, including those computed in the bracketing line search. All methods were terminated after 10^4 function evaluations. Each function evaluation $f(y)$ makes a call to the MATLAB function `eig` to compute all the eigenvalues of $C - \mathcal{A}^T y$. The vertical axis shows the relative error $|(f - f^*)/f^*|$, where we used the SDPT3 solver in CVX to obtain the optimal solution f^* with accuracy 10^{-14} . As in the experiment on the nonsmooth Les Houches problem reported in Figure 1, scaled L-BFGS-1 (magenta dots) breaks down early. However, unlike in that experiment, here unscaled L-BFGS-1 (cyan dots) also breaks down early, and as a result the subgradient method (green dots) obtains a better answer, though not nearly as good as full BFGS (black dots).

It's also of interest to examine the multiplicity of the eigenvalues of $C - \mathcal{A}^T y$ at the optimal solution y^* and its computed approximations. From SDPT3, we know that the optimal multiplicity for this problem is 5. Table 1 shows the top 6 eigenvalues of the final answer found by each method. Besides SDPT3, only BFGS and the subgradient method are able to determine the correct optimal multiplicity. BFGS finds a solution with 11 correct digits, while the subgradient method obtains 3 correct digits. Both variants of L-BFGS-1 converge to answers with multiplicity 4. Although the multiplicity is wrong, unscaled L-BFGS-1 gets a better answer than its scaled counterpart, with 2 correct digits, but at the cost of requiring many more function evaluations.

Next, we repeat this experiment on the same problem, increasing the number of L-BFGS updates from $m = 1$ to $m = 20$. See Figure 6 as well as Table 2 which presents the top 6 eigenvalues for the final answer obtained by scaled and unscaled L-BFGS-20. Both methods find the right multiplicity, with scaled L-BFGS-20 obtaining 3 correct digits and unscaled L-BFGS-20 obtaining 4 correct digits.

In summary, we observe that for the Max Eigenvalue problem, unlike the Les Houches problem, increasing m from 1 to 20 does not result in scaled L-BFGS doing better than unscaled L-BFGS.

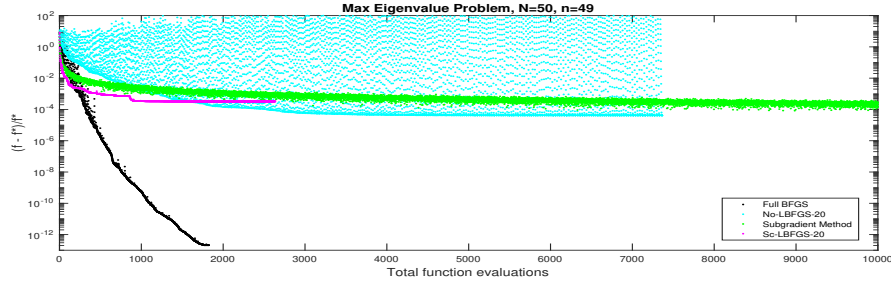


Figure 6: Comparing BFGS, L-BFGS-20 with and without scaling and subgradient method on a randomly generated Max Eigenvalue problem (18) with $N = 50$ and $n = 49$.

Sc L-BFGS-20	No L-BFGS-20
7.82959659952176	7.82735384547039
7.82959659952176	7.82735380191247
7.82959659952174	7.82735377961470
7.82959659952166	7.82735377236995
7.82959659950328	7.82735372682267
7.62982269438813	7.69181664817458

Table 2: Top 6 eigenvalues of $C - \mathcal{A}^T y$ for Max Eigenvalue problem (18) where y is computed by scaled and unscaled L-BFGS-20 for the same problem reported in Table 1. The optimal multiplicity is 5.

3.4 Smoothed Max Eigenvalue Problem

Consider now Nesterov smoothing of the Max Eigenvalue function (18) [d'A08]

$$f_\mu(y) = \mu \log \sum_{i=1}^N \exp(\lambda_i(C - \mathcal{A}^T y)/\mu) - \mu \log N, \quad (19)$$

where $\lambda_1(W) \geq \lambda_2(W) \geq \dots \geq \lambda_N(W)$ denote the ordered eigenvalues of a symmetric matrix $W \in S^N$. Thus, λ_1 is equivalent to λ_{\max} . Unlike the Les Houches problem, where the nonsmooth optimal value is equal to the smoothed optimal value, that is $f^* = f_\mu^* = 0$, for any μ as $\mu \rightarrow 0$, the same statement is not true for the Max Eigenvalue problem. The smoothed Max Eigenvalue problem requires a complete eigendecomposition in order to obtain every eigenvalue for a given matrix, and since CVX does not allow such functions but only those that it knows to be convex such as the maximum eigenvalue function, we could not compute the optimal value f_μ^* from CVX. Instead, we use full BFGS with the maximum number of iterations set to 10^5 to minimize (19) to high accuracy: we denote this computed value by f_μ^B .

In Figure 7 we report on an experiment using the smoothed version of the same instance of the randomly generated Max Eigenvalue problem as earlier with $N = 50$ and $n = 49$, using L-BFGS-1 to minimize (19). The left panel shows the final value computed by scaled (magenta asterisks) and unscaled (cyan squares) L-BFGS-1 shifted by f_μ^B (the answer found by full BFGS), as a function of the smoothing parameter μ , in log-log scale. The right panel shows the number of iterations as a function of μ , also in log-log scale. The maximum number of iterations is 10^5 .

In the left panel we see that for $\mu = 1$ down to $\mu = 10^{-4}$ both methods yield about the same accuracy as each other, but that this deteriorates as μ decreases. Scaled L-BFGS-1 continues to obtain a reasonable approximation to the presumed accurate solution f_μ^B for μ down to 10^{-9} , although this accuracy continues to decrease as μ is reduced. Looking at the right panel, we see that starting with $\mu = 10^{-10}$ scaled L-BFGS-1 hits the maximum

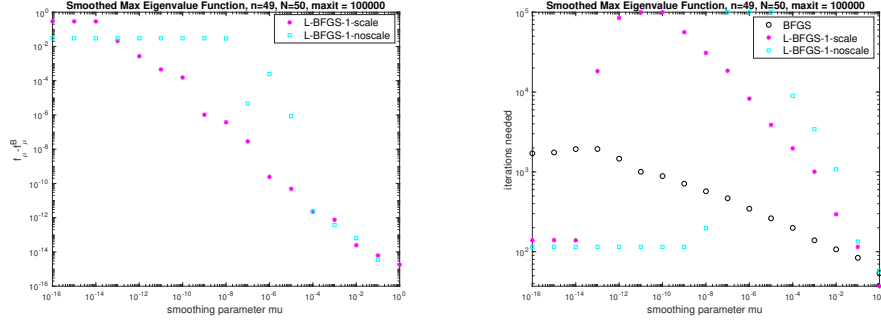


Figure 7: Comparing L-BFGS-1 with and without scaling on the smoothed Max Eigenvalue problem (19) for $N = 50$ and $n = 49$. The left panel shows the final function value, shifted by f_μ^B , the optimal value computed by BFGS, and the right panel shows the iteration count, both as a function of the smoothing parameter μ . The maximum number of iterations is set to 10^5 .

SDPT3	BFGS	Sc L-BFGS-1	No L-BFGS-1
7.82702970305352	7.82702976093363	7.82702978971152	7.82703432112405
7.82702970305349	7.82702968960443	7.82702971836333	7.82703424950540
7.82702970305348	7.82702966768912	7.82702969644540	7.82703422776180
7.82702970305346	7.82702963829977	7.82702966707913	7.82703419808905
7.82702970305334	7.82702954704101	7.82702957585856	7.82703410710019
7.70350538019538	7.70350539089203	7.70374015675041	7.69886385782543

Table 3: Top 6 eigenvalues of $C - \mathcal{A}^T y$ for Max Eigenvalue problem where y is computed by applying BFGS, scaled L-BFGS-1 and unscaled L-BFGS-1 to f_μ with $\mu = 10^{-7}$, for the same instance of the randomly generated Max Eigenvalue problem as in Table 1. The optimal multiplicity is 5. The first column gives the top 6 eigenvalues of the solution to the original nonsmooth problem.

iteration limit and starting with 10^{-12} it breaks down before reaching the maximum iteration limit. In contrast, unscaled L-BFGS-1 hits the maximum iteration limit for $\mu = 10^{-5}$ and breaks down for $\mu \leq 10^{-8}$.

Table 3 shows the top 6 eigenvalues of the final answer found by BFGS and L-BFGS-1 for the smoothed Max Eigenvalue problem with $\mu = 10^{-7}$. We also repeat the optimal top 6 eigenvalues of the minimizer of the original nonsmooth function f found by SDPT3 for the sake of comparison. Note that the result computed by applying scaled L-BFGS-1 to the smoothed problem agrees with the nonsmooth optimal value f^* to 8 digits, compared to 0 digits when applied directly to the nonsmooth problem (Table 1).

We repeated this experiment with $m = 20$, reported in Figure 8. In the left panel, we observe that the loss of accuracy in L-BFGS as a function of μ is less pronounced with 20 updates. Roughly speaking, overall the error decreases by a factor of 10^{-2} . In the right panel we see that neither scaled nor unscaled L-BFGS-20 reaches the maximum iteration limit, but that both methods break down for sufficiently small μ . The top eigenvalues produced by L-BFGS-20 for $\mu = 10^{-7}$ are shown in Table 4.

Comparing the final computed maximum eigenvalue in Tables 1, 2, 3, and 4, note that full BFGS obtains a more accurate solution when applied to the original nonsmooth problem than it does when applied to the smoothed approximation using $\mu = 10^{-7}$, while the opposite is true for scaled L-BFGS-1 and scaled L-BFGS-20. In summary, as with the Les Houches problem, it is much more effective to apply L-BFGS to the smoothed max eigenvalue problem than directly to the nonsmooth problem. As earlier, this is in sharp contrast to the behavior of full BFGS.

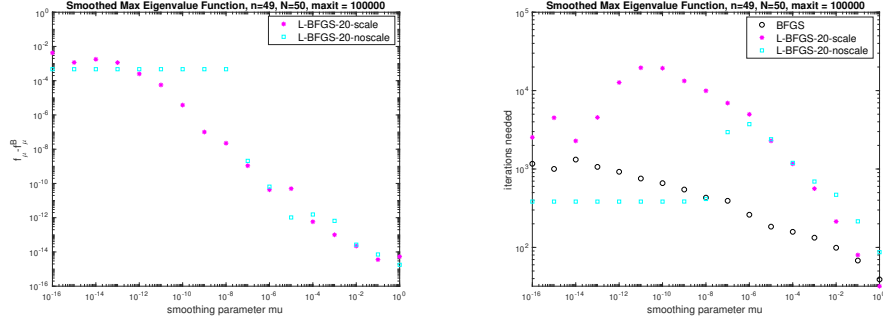


Figure 8: Comparing L-BFGS-20 with and without scaling on the smoothed Max Eigenvalue problem (19) for $N = 50$ and $n = 49$. The left panel shows the final function value, shifted by f_μ^B , the optimal value computed by BFGS, and the right panel shows the iteration count, both as a function of the smoothing parameter μ . The maximum number of iterations is set to 10^5 .

Sc L-BFGS-20	No L-BFGS-20
7.82702976201688	7.82702976326908
7.82702969067290	7.82702969273000
7.82702966877706	7.82702966880138
7.82702963938214	7.82702964120678
7.82702954813253	7.82702954355970
7.70348252862186	7.70344821498698

Table 4: Top 6 eigenvalues of $C - \mathcal{A}^T y$ for Max Eigenvalue problem where y is computed by applying BFGS, scaled L-BFGS-20 and unscaled L-BFGS-20 to f_μ with $\mu = 10^{-7}$, for the same instance of the randomly generated Max Eigenvalue problem as in Table 1. The optimal multiplicity is 5.

3.5 Semidefinite Programming

Consider the following primal and dual semidefinite programs (SDP) in standard form [HR00, HOR14]

$$\max_{X \in S^N} \langle C, X \rangle \quad (20)$$

$$\text{subject to } \mathcal{A}X = b \quad \text{and} \quad X \in S_+^N,$$

$$\min_{y \in \mathbb{R}^n} b^T y \quad (21)$$

$$\text{subject to } Z = \mathcal{A}^T y - C \quad \text{and} \quad Z \in S_+^N,$$

where $b \in \mathbb{R}^n$, $C \in S^N$ and $\mathcal{A} : S^N \rightarrow \mathbb{R}^n$ is a linear operator as defined in (16) and (17). Here $S_+^N \subseteq S^N$ denotes the cone of positive semidefinite $N \times N$ matrices. Let us assume that strong duality holds, so that the optimal primal and dual values are the same, and that the optimal values are attained. It follows that if X^* is an optimal solution to the primal problem (20) and Z^* is an optimal solution to the dual problem (21), we have $X^* Z^* = 0$. Further assume that X^* is nonzero, and consequently, Z^* has at least one eigenvalue equal to zero. Then the dual problem (21) is equivalent to the following unconstrained eigenvalue optimization problem

$$\min_{y \in \mathbb{R}^n} f(y), \quad (22)$$

with the exact penalty dual function [DYC⁺19]

$$f(y) = b^T y + \alpha \max\{\lambda_{\max}(C - \mathcal{A}^T y), 0\}, \quad (23)$$

for sufficiently large α , where λ_{\max} denotes maximum eigenvalue as earlier. Note that this exact penalty function differs from the eigenvalue optimization formulation in [HR00], namely

$$b^T y + \alpha \lambda_{\max}(C - \mathcal{A}^T y)$$

which does not include the $\max\{\cdot, 0\}$ operator. In that formulation, to give a correct equivalence α must be exactly equal to a critical value, as opposed to greater than or equal to this value. For the SDP problems we consider in the following subsections we already know valid lower bounds for α . Note that at an optimal solution y^* the maximum eigenvalue of $-Z^* = C - \mathcal{A}^T y^*$ is zero, often with multiplicity greater than one, and hence f is nonsmooth at y^* .

3.6 Max Cut Problem

Our first example of semidefinite programming (SDP) arises from the Max Cut problem. This subsection and a subsequent one on the Matrix Completion problem were motivated by the recent paper [DYC⁺19] and the observation made there that the first-order algorithms they used to minimize the penalized dual function (23) arising from Max Cut and Matrix Completion SDP relaxations were slow. Here we compare full BFGS, scaled and unscaled L-BFGS and the subgradient method (again with $t_k = 1/k$) on penalized dual functions arising from Goemans-Williamson SDP relaxations of the Max Cut problem. We note that one of the key ideas in [DYC⁺19] is that, when the primal SDP optimal solution X^* has rank much less than N , an accurate estimate of the optimal value of the SDP obtained from minimizing the penalized dual function allows the use of a novel method for obtaining efficient low-rank solutions to the *primal* SDP even when N is large.

The primal Max Cut SDP relaxation and its dual are [GW95]:

$$\max_X \quad \frac{1}{4} \langle L, X \rangle \tag{24}$$

$$\text{subject to} \quad \text{diag}(X) = \mathbf{1} \quad \text{and} \quad X \in S_+^N,$$

$$\min_{y \in \mathbb{R}^n} \quad \mathbf{1}^T y \tag{25}$$

$$\text{subject to} \quad Z = \text{Diag}(y) - \frac{1}{4}L \quad \text{and} \quad Z \in S_+^N,$$

where L is the Laplacian matrix of a given undirected graph, $\text{diag}(\cdot)$ maps the diagonal of a matrix to a vector, and $\text{Diag}(\cdot)$ maps a vector to a diagonal matrix. Note that these are instances of the primal and dual SDP introduced in (20) and (21), respectively. By definition for the SDP Max Cut problem we have $n = N$. The exact penalty dual function (23) for the Max Cut SDP relaxation is

$$f(y) = \mathbf{1}^T y + \alpha \max\{\lambda_{\max}(L - \text{Diag}(y)), 0\}. \tag{26}$$

Due to the constant trace property of the primal Max Cut SDP (24), the trace (nuclear) norm of the primal optimal solution is known to be N , and hence any solution y^* to the penalized dual max cut problem (26) with $\alpha \geq N$ is also a solution to the dual SDP (25) and vice versa [DYC⁺19, Lem. 6.1].

We picked graph $G1$ from the Gset group in the sparse matrix collection [Gse14] for the following experiment. $G1$ is an unweighted graph with $N = 800$ vertices and the adjacency matrix is a sparse symmetric matrix with 38352 nonzero entries (all equal to 1). Since N is relatively small we can apply the SDPT3 solver via CVX [GB08] to the primal SDP (24), obtaining the optimal primal and dual value $f^* = 12083.19765$. The rank r^* of the

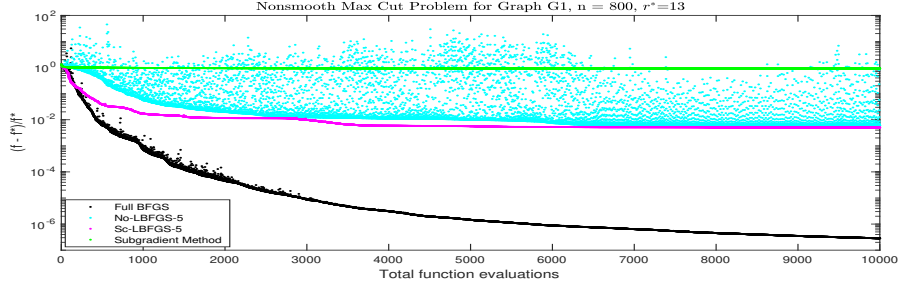


Figure 9: Comparing BFGS, L-BFGS-5 with and without scaling and the subgradient method on the penalized dual Max Cut problem (26).

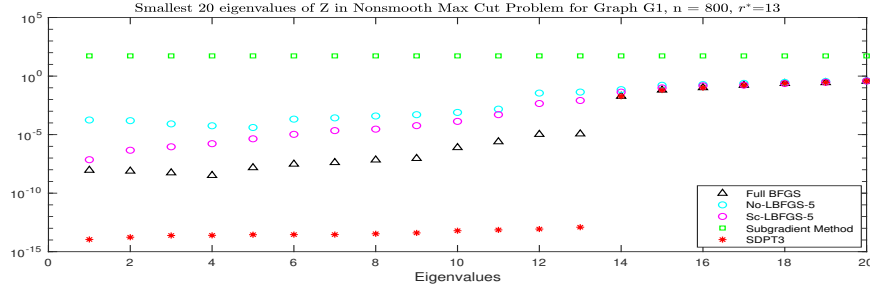


Figure 10: Comparing smallest 20 eigenvalues of the dual slack matrix Z obtained by BFGS, L-BFGS-5 with and without scaling and the subgradient method on the penalized dual Max Cut problem (26) for the G_1 graph with $n = 800$. The nullity of the optimal dual slack matrix Z^* is 13. The smallest 20 eigenvalues obtained from SDPT3 are shown as well. The lack of monotonicity at the left end of some of the plots occurs because we actually plotted the absolute values of the ordered largest eigenvalues of $-Z$, and some of these eigenvalues are positive, either because of rounding errors or insufficient accuracy in the optimization.

optimal primal solution X^* is 13, and strict complementarity holds, so the nullity of the dual solution Z^* is also 13.

In Figure 9, we compare the performance of full BFGS, L-BFGS-5 with and without scaling, and the subgradient method (with $t_k = 1/k$) to minimize the penalized dual function (26) with $\alpha = 2N = 1600$. We display all the function values that were computed, with the maximum number of function evaluations set to 10^4 . The vertical axis shows the relative error $(f - f^*)/f^*$. In contrast to the two experiments presented in Figure 5 and 6 for the nonsmooth Max Eigenvalue problem, the results of this experiment are in favor of L-BFGS when compared to the subgradient method. Both scaled L-BFGS-5 (magenta dots) and unscaled L-BFGS-5 (cyan dots) reduce the relative error down to below 10^{-2} while the subgradient method (blue dots) reduces the relative error to about 10^0 . Full BFGS (black dots) reduces the relative error to below 10^{-6} .

In Figure 10, we show the *negative* of the top 20 eigenvalues of the final negative dual slack matrix $-Z$ (equivalently, the smallest 20 eigenvalues of Z) obtained by the four methods, along with values obtained by SDPT3. It is interesting to note that full BFGS approximates the eigenvalues well, in the sense that it clearly separates the first 13 approximately zero eigenvalues from the approximations to the nonzero eigenvalues, although not as decisively as SDPT3. However, the final eigenvalues obtained by L-BFGS-5 are not clearly separated.

Next we increase the number of L-BFGS updates from $m = 5$ to $m = 20$, showing the results in Figures 11 and 12. The results for BFGS and the subgradient method are shown again for comparison. We see that scaled L-BFGS-20 now reduces the relative error down to 10^{-4} and unscaled to about 10^{-3} , compared to about 10^{-2} for scaled and unscaled L-

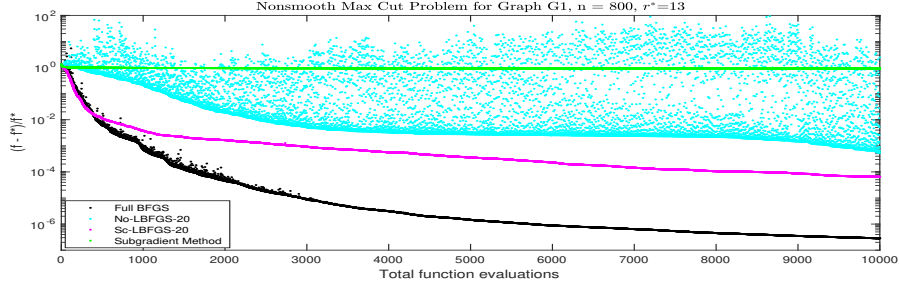


Figure 11: Comparing BFGS, L-BFGS-20 with and without scaling and the subgradient method on the penalized dual Max Cut problem (26).

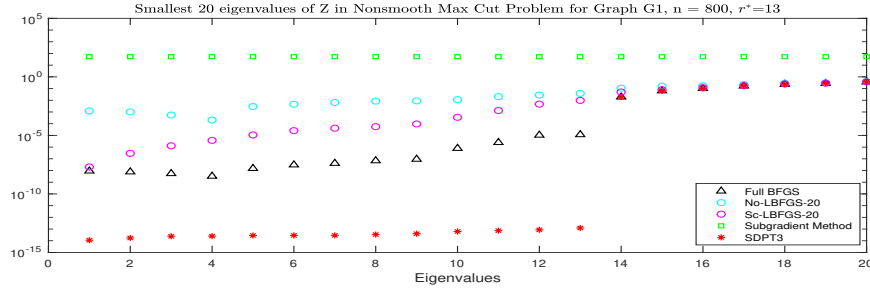


Figure 12: Comparing smallest 20 eigenvalues of the dual slack matrix Z obtained by BFGS, L-BFGS-5 with and without scaling and the subgradient method on the penalized dual Max Cut problem (26) for $G1$ graph with $n = 800$. The nullity of the optimal dual slack matrix Z^* is 13. The smallest 20 eigenvalues obtained from SDPT3 are shown as well. See legend of Figure 10 regarding eigenvalue monotonicity.

BFGS-5. However, the eigenvalue plot is similar to the corresponding plot for L-BFGS-5: neither variant is able to discover that the nullity of the optimal dual slack matrix Z^* is 13.

3.7 Smoothed Max Cut Problem

Consider now Nesterov smoothing of the penalized dual Max Cut problem (26):

$$f_\mu(y) = \mathbf{1}^T y + \alpha \mu \log \left(1 + \sum_{i=1}^n \exp(\lambda_i(L - \text{Diag}(y))/\mu) \right) - \alpha \mu \log(n+1). \quad (27)$$

Note the presence of the term “1” which does not appear in (19): this reflects the presence of the $\max\{\cdot, 0\}$ operator in the penalty function (23).

Figure 13 shows the results of applying BFGS, scaled L-BFGS-5 and scaled L-BFGS-20 to minimize (27) with $\mu = 10^{-7}$. We do not include the unscaled L-BFGS variants because it seems clear that they offer no advantage on smooth problems. Interestingly, and in contrast to our other experiments with smoothed nonsmooth functions, scaled L-BFGS-20 does better than full BFGS during its first several thousand function evaluations, but eventually it is overtaken by full BFGS. Scaled L-BFGS-5 does relatively poorly.

Table 5 shows the final answers we obtained from full BFGS, scaled L-BFGS-20 and scaled L-BFGS-5 on the smoothed and nonsmooth Max Cut problem. The optimal SDP value f^* is also shown. All three of full BFGS, scaled L-BFGS-5 and scaled L-BFGS-20 obtain better approximations to f^* when applied directly to the nonsmooth problem than when applied to the smoothed problem with $\mu = 10^{-7}$. This is in contrast to what we observed for the Les Houches problem and the Max Eigenvalue problem, where this was true only for full BFGS.

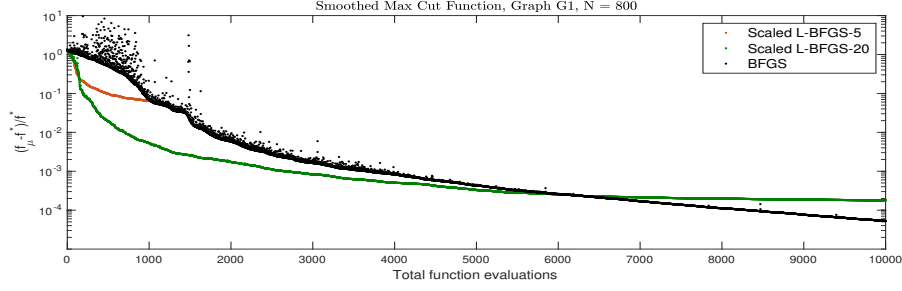


Figure 13: Comparing scaled L-BFGS-5, scaled L-BFGS-20 and BFGS on Smoothed Max Cut problem (27) for the same problem as in the Nonsmooth Max Cut problem. The smoothing parameter is $\mu = 10^{-7}$.

SDP-optimal	12083.19765454945
BFGS-smooth	12083.83341694415
BFGS-nonsmooth	12083.20108505506
L-BFGS-20-smooth	12085.35533081401
L-BFGS-20-nonsmooth	12083.97779371002
L-BFGS-5-smooth	12848.47893036591
L-BFGS-5-nonsmooth	12143.81352524515

Table 5: Final objective value we obtained from full BFGS, scaled L-BFGS-20 and scaled L-BFGS-5 on the smoothed and nonsmooth Max Cut problem presented in Figures 13, 11 and 9 respectively. The optimal SDP value f^* is also shown for comparison.

It would be interesting to investigate whether L-BFGS might be useful in the solution of large-scale Max Cut problems or other SDPs. A first step in this direction appears in [Asl20, Sec. 4.2.5], utilizing a variant of the smoothed exact penalty dual function (27) that includes only the largest eigenvalues in the smoothing, since these are the ones that dominate (27). This allows the use of MATLAB's `eigs` to compute only the largest few eigenvalues of $C - \mathcal{A}^T y$ via the Lanczos method.

3.8 Matrix Completion Problem

The Matrix Completion problem is as follows. Suppose $\mathcal{X} \in \mathbb{R}^{N_1 \times N_2}$ denotes a low-rank matrix for which we only have access to some of its entries and would like to recover entirely by minimizing the rank over all matrices whose entries agree with the known values. This rank minimization problem is NP-hard, so we relax it by minimizing a well-known convex surrogate for the rank: the nuclear norm (sum of all the singular values). Let Ω be the set of pairs (i, j) for which \mathcal{X}_{ij} is known. Then the nuclear norm minimization problem can be expressed as the following SDP [RFP10]

$$\begin{aligned}
 & \max_{X \in S^{(N_1+N_2)}} && -\text{Tr}(W_1) - \text{Tr}(W_2) \\
 & \text{subject to} && U_{ij} = \mathcal{X}_{ij}, (i, j) \in \Omega, \\
 & && X = \begin{bmatrix} W_1 & U \\ U^T & W_2 \end{bmatrix} \in S_+^{(N_1+N_2)}.
 \end{aligned} \tag{28}$$

We write the primal problem in the *max* form in order to be consistent with the SDP form (20), with $N = N_1 + N_2$. Define the constraint $U_{ij} = \mathcal{X}_{ij}$ for $(i, j) \in \Omega$ in linear operator form $\mathcal{B}(U) = b$, where $b \in \mathbb{R}^n$ with $n = |\Omega|$ is the vector consisting of the known entries of \mathcal{X} in some prescribed order and $\mathcal{B} : \mathbb{R}^{N_1 \times N_2} \rightarrow \mathbb{R}^n$, with \mathcal{B}^T its adjoint operator. The dual

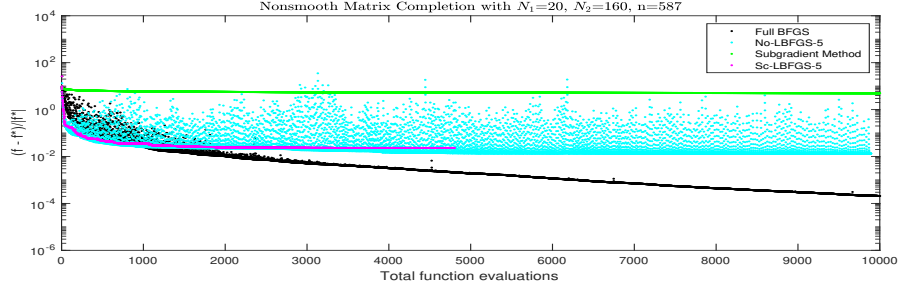


Figure 14: Comparing BFGS, LBFGS-5 with and without scaling and the subgradient method on the penalized dual Matrix Completion problem (30).

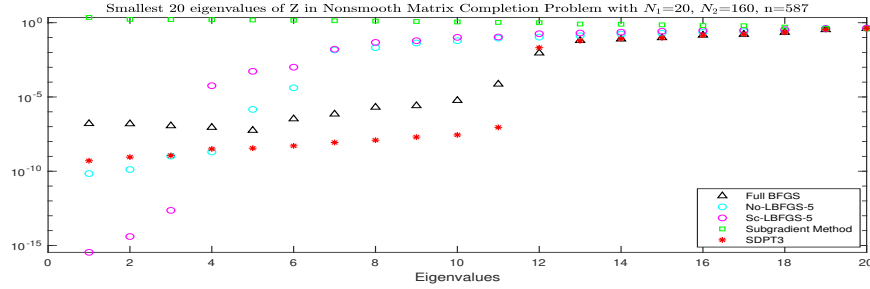


Figure 15: Comparing smallest 20 eigenvalues of the dual slack matrix Z obtained by BFGS, L-BFGS-5 with and without scaling and the subgradient method on the penalized dual Matrix Completion problem (30). See legend of Figure 10 regarding eigenvalue monotonicity.

SDP is

$$\begin{aligned} \min_{y \in \mathbb{R}^n} \quad & b^T y \\ \text{subject to} \quad & Z = \begin{bmatrix} I_{N_1} & \mathcal{B}^T(y) \\ (\mathcal{B}^T(y))^T & I_{N_2} \end{bmatrix}, \quad Z \in S_+^{(N_1+N_2)}. \end{aligned} \quad (29)$$

The exact penalty dual function (23) for the Matrix Completion problem is then

$$f(y) = b^T y + \alpha \max \left\{ \lambda_{\max} \left(- \begin{bmatrix} I_{N_1} & \mathcal{B}^T(y) \\ (\mathcal{B}^T(y))^T & I_{N_2} \end{bmatrix} \right), 0 \right\}. \quad (30)$$

For the experiment in this part, we generated a low-rank random matrix \mathcal{X} of size $N_1 = 20$ by $N_2 = 160$ with rank 3. We then selected the ordered pairs in Ω randomly with the probability of each (i, j) being included set to 0.2. For this problem instance we got $|\Omega| = n = 587$. We then applied the various methods to minimize (30) with $\alpha = 2\|X^*\|_* = 2\text{Tr}(X^*) = -2f^* = 3.0796$, where $f^* = -1.5398$ and $X^* \in S_+^{N_1+N_2}$ were obtained from solving the SDP (29) via SDPT3. Note that the optimal value is negative because of the minus sign in the \max formulation of the primal SDP.

Figure 14 shows the performance of full BFGS, L-BFGS-5 with and without scaling, and the subgradient method (with $t_k = 1/k$). The vertical axis shows the relative error $(f - f^*)/|f^*|$. The maximum number of function evaluations is set to 10^4 . As is evident from the plot, both variants of L-BFGS-5 outperform the subgradient method, even though scaled L-BFGS-5 quits early before 5000 evaluations and unscaled L-BFGS-5 just before 10000 evaluations.

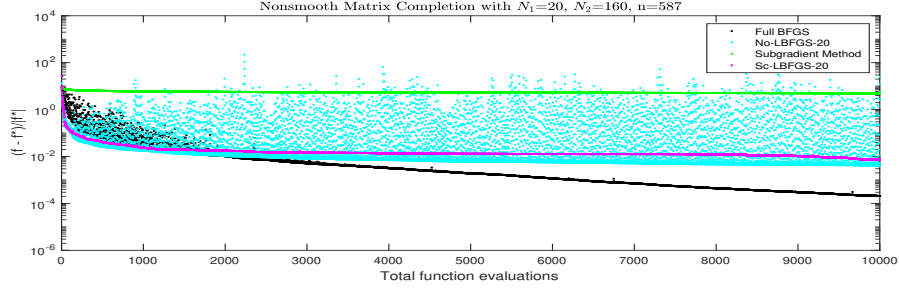


Figure 16: Comparing BFGS, LBFGS-20 with and without scaling and subgradient method on the penalized dual Matrix Completion problem (30).

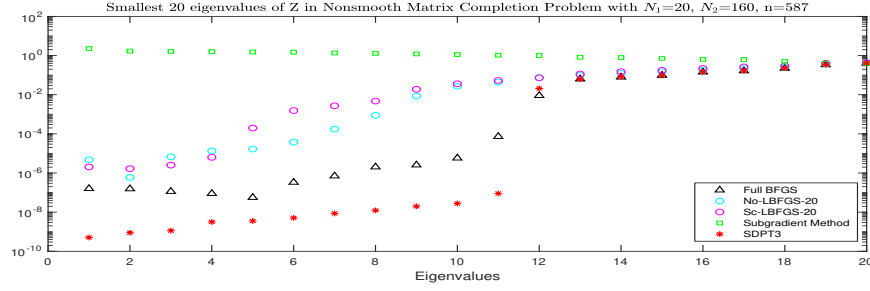


Figure 17: Comparing smallest 20 eigenvalues of the dual slack matrix Z obtained by BFGS, L-BFGS-20 with and without scaling and the subgradient method on the penalized dual Matrix Completion problem (30). See legend of Figure 10 regarding eigenvalue monotonicity.

Figure 15 presents the negative of the top 20 eigenvalues of the final negative dual slack matrix $-Z$, or equivalently, the 20 smallest eigenvalues of Z , obtained by the four methods, along with values obtained by SDPT3. As before, BFGS is able to separate the zero and nonzero eigenvalues of Z^* , agreeing with SDPT3 that the nullity of Z^* is effectively 11. Note that this is larger than 3, the rank of the original matrix \mathcal{X} , implying that 20% was not enough observations to reconstruct \mathcal{X} . It is interesting that the eigenvalues of the solution found by scaled L-BFGS-5 do suggest a nullity of 3, but this may just be a coincidence.

In the experiment reported in Figures 16 and 17, we increase m to 20 and again we compare the relative error and the smallest 20 eigenvalues of the dual slack matrix, respectively. In both plots the result from BFGS and the subgradient method are repeated for comparison. In Figure 16, neither L-BFGS-20 method quits early this time; the unscaled variant gets a slightly lower answer. The eigenvalues shown for L-BFGS-20 in Figure 17 do not suggest any conclusion about the nullity of Z^* .

3.9 Smoothed Matrix Completion Problem

Nesterov smoothing of the penalized dual Matrix Completion problem (30) gives

$$\begin{aligned}
 f_\mu(y) &= b^T y \\
 &+ \alpha \mu \log \left(1 + \sum_{i=1}^{N_1+N_2} \exp \left(\lambda_i \left(- \begin{bmatrix} I_{N_1} & \mathcal{B}^T(y) \\ (\mathcal{B}^T(y))^T & I_{N_2} \end{bmatrix} / \mu \right) \right) \right) \\
 &- \alpha \mu \log(N_1 + N_2 + 1).
 \end{aligned} \tag{31}$$

In Figure 18, we show the result of applying full BFGS, scaled L-BFGS-5 and scaled L-BFGS-20 to (31) with $\mu = 10^{-7}$. The underlying reference matrix, \mathcal{X} , is the same matrix

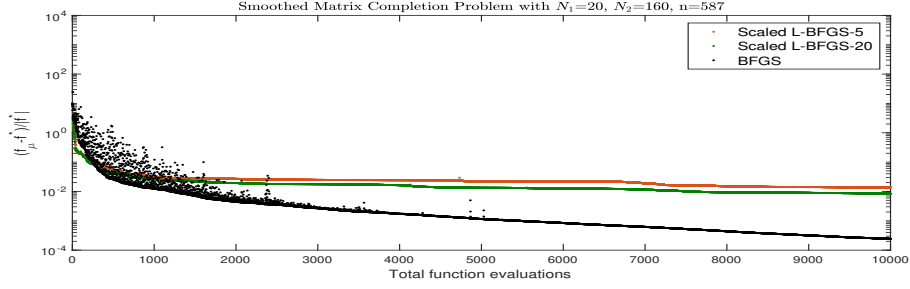


Figure 18: Comparing scaled L-BFGS-5, scaled L-BFGS-20 and BFGS on Smoothed Matrix Completion problem (31) for the same problem as in the Nonsmooth Matrix Completion problem. The smoothing parameter is $\mu = 10^{-7}$.

SDP-optimal	-1.5397855575175
BFGS-smooth	-1.53941270679104
BFGS-nonsmooth	-1.53946718487809
L-BFGS-20-smooth	-1.52686081626082
L-BFGS-20-nonsmooth	-1.52864965852708
L-BFGS-5-smooth	-1.51868588317043
L-BFGS-5-nonsmooth	-1.50422184883968

Table 6: Final objective value we obtained from full BFGS, scaled L-BFGS-20 and scaled L-BFGS-5 on the smoothed and nonsmooth Matrix Completion problem presented in Figures 18, 16 and 14 respectively. The optimal SDP value f^* is also shown for comparison.

as in the nonsmooth experiment in Figure 14, with $N_1 = 20$, $N_2 = 160$, $f^* = -1.5398$. We set $\alpha = 3.0796$ as before. The maximum number of function evaluations is set to 10^4 .

We see in Figure 18 that although the relative error obtained by scaled L-BFGS-20 and full BFGS are about the same as when applied to the nonsmooth function, scaled L-BFGS-5 gets a lower error when it is applied to the smoothed function, and more importantly, it does not break down early.

Table 6 shows the final answers we obtained from full BFGS, scaled L-BFGS-20 and scaled L-BFGS-5 on the smoothed and nonsmooth Matrix Completion problem. The optimal SDP value f^* is also shown for comparison.

4 Concluding Remarks

In §2, we presented theoretical results showing that L-BFGS may converge to non-optimal points when applied to a simple class of nonsmooth functions. In §3, we investigated whether the same phenomenon holds in practice. We found that when applied to a nonsmooth function directly, L-BFGS, especially its scaled variant, often breaks down with a poor approximation to an optimal solution, in sharp contrast to full BFGS. Unscaled L-BFGS is less prone to breakdown but conducts far more function evaluations per iteration than scaled L-BFGS does, and thus it is slow. Nonetheless, it is often the case that both variants obtain better results than the provably convergent, but slow, subgradient method.

On the other hand, when applied to a smooth approximation of a nonsmooth function, scaled L-BFGS invariably performs better than unscaled L-BFGS, often obtaining good results even when the problem is quite ill-conditioned. In particular, scaled L-BFGS may be a reasonable approach to finding approximate minimizers of smoothed exact penalty dual functions arising in large-scale semidefinite programs, although further investigation is needed to investigate the practicality of this approach. Minimization of the SDP exact

penalty dual function is a key component of a recently proposed method for solving large-scale SDPs with low-rank primal solutions [DYC⁺19].

Most importantly, we find that although L-BFGS is often a reliable method for minimizing ill-conditioned smooth problems, it frequently fails when the condition number is so large that the function is effectively nonsmooth. This behavior is in sharp contrast to the behavior of full BFGS, which is consistently reliable for nonsmooth optimization problems. We arrive at the conclusion that, for large-scale nonsmooth optimization problems for which full BFGS and other methods are not practical, it is often better to apply L-BFGS to a smoothed variant of a nonsmooth problem than to apply it directly to the nonsmooth problem.

Acknowledgments. The work of the first author was conducted as part of her Ph.D. studies at the Courant Institute of Mathematical Studies, New York University. Many thanks to Margaret H. Wright for providing her financial support through a grant from the Simons Foundation (417314,MHW).

References

- [AO20a] Azam Asl and Michael L. Overton. Analysis of limited-memory BFGS on a class of nonsmooth convex functions. *IMA Journal of Numerical Analysis*, 01 2020. drz052.
- [AO20b] Azam Asl and Michael L. Overton. Analysis of the gradient method with an Armijo-Wolfe line search on a class of non-smooth convex functions. *Optimization Methods and Software*, 35(2):223–242, 2020.
- [Asl20] Azam Asl. *Behavior of the Limited Memory BFGS Method on Nonsmooth Optimization Problems in Theory and Practice*. PhD thesis, New York University, 2020. https://cs.nyu.edu/media/publications/asl.thesis_final_UtpoLsu.pdf.
- [BB88] Jonathan Barzilai and Jonathan M. Borwein. Two-point step size gradient methods. *IMA Journal of Numerical Analysis*, 8(1):141–148, 1988.
- [BCL⁺20] James V. Burke, Frank E. Curtis, Adrian S. Lewis, Michael L. Overton, and Lucas E. A. Simões. Gradient sampling methods for nonsmooth optimization. In Adil M. Bagirov, Manlio Gaudioso, Napsu Karmitsa, Marko M. Mäkelä, and Sona Taheri, editors, *Numerical Nonsmooth Optimization: State of the Art Algorithms*, pages 201–225. Springer International Publishing, Cham, 2020.
- [BGK⁺20] Adil M. Bagirov, Manlio Gaudioso, Napsu Karmitsa, Marko M. Mäkelä, and Sona Taheri. *Numerical Nonsmooth Optimization: State of the Art Algorithms*. Springer International Publishing, Cham, 2020.
- [BKM14] Adil Bagirov, Napsu Karmitsa, and Marko M. Mäkelä. *Introduction to Nonsmooth Optimization*. Springer, Cham, 2014. Theory, Practice and Software.
- [BLO05] James V. Burke, Adrian S. Lewis, and Michael L. Overton. A robust gradient sampling algorithm for nonsmooth, nonconvex optimization. *SIAM J. Optim.*, 15(3):751–779, 2005.
- [BV04] Stephen Boyd and Lieven Vandenbergh. *Convex Optimization*. Cambridge University Press, March 2004.
- [CL20] Frank E. Curtis and Minhan Li. Gradient sampling methods with inexact subproblem solutions and gradient aggregation, 2020. arXiv:2005.07822.

- [CMO17] Frank E. Curtis, Tim Mitchell, and Michael L. Overton. A BFGS-SQP method for nonsmooth, nonconvex, constrained optimization and its evaluation using relative minimization profiles. *Optimization Methods and Software*, 32(1):148–181, 2017.
- [CO12] Frank E. Curtis and Michael L. Overton. A sequential quadratic programming algorithm for nonconvex, nonsmooth constrained optimization. *SIAM J. Optim.*, 22(2):474–500, 2012.
- [d’A08] Alexandre d’Aspremont. Smooth optimization with approximate gradient. *SIAM J. Optim.*, 19(3):1171–1183, 2008.
- [DYC⁺19] Lijun Ding, Alp Yurtsever, Volkan Cevher, Joel A. Tropp, and Madeleine Udell. An optimal-storage approach to semidefinite programming using approximate complementarity, February 2019. arXiv:1902.03373.
- [GB08] Michael Grant and Stephen Boyd. Graph implementations for non-smooth convex programs. In V. Blondel, S. Boyd, and H. Kimura, editors, *Recent Advances in Learning and Control*, Lecture Notes in Control and Information Sciences, pages 95–110. Springer-Verlag Limited, 2008. http://stanford.edu/~boyd/graph_dcp.html.
- [GB14] Michael Grant and Stephen Boyd. CVX: Matlab software for disciplined convex programming, version 2.1. <http://cvxr.com/cvx>, March 2014.
- [GL18] J. Guo and A. Lewis. Nonsmooth variants of Powell’s BFGS convergence theorem. *SIAM Journal on Optimization*, 28(2):1301–1311, 2018.
- [Gse14] The University of Florida sparse matrix collection: Gset group. <http://www.cise.ufl.edu/research/sparse/matrices/Gset/index.html>, March 2014.
- [GW95] Michel X. Goemans and David P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *J. ACM*, 42(6):1115–1145, November 1995.
- [HOR14] C. Helmberg, M.L. Overton, and F. Rendl. The spectral bundle method with second-order information. *Optimization Methods and Software*, 29(4):855–876, 2014.
- [HR00] C. Helmberg and F. Rendl. A spectral bundle method for semidefinite programming. *SIAM Journal on Optimization*, 10(3):673–696, 2000.
- [Kiw85] Krzysztof C. Kiwiel. *Methods of descent for nondifferentiable optimization*, volume 1133 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 1985.
- [Lem75] C. Lemaréchal. An extension of Davidon methods to non differentiable problems. *Math. Programming Stud.*, (3):95–109, 1975.
- [LN89] Dong C. Liu and Jorge Nocedal. On the limited memory BFGS method for large scale optimization. *Math. Programming*, 45(3, (Ser. B)):503–528, 1989.
- [LO13] Adrian S. Lewis and Michael L. Overton. Nonsmooth optimization via quasi-Newton methods. *Math. Program.*, 141(1-2, Ser. A):135–163, 2013.
- [LZ15] A. S. Lewis and S. Zhang. Nonsmoothness and a variable metric method. *J. Optim. Theory Appl.*, 165(1):151–171, 2015.

- [MRT12] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. Adaptive Computation and Machine Learning. MIT Press, Cambridge, MA, 2012.
- [NB01] Angelia Nedić and Dimitri P. Bertsekas. Incremental subgradient methods for nondifferentiable optimization. *SIAM J. Optim.*, 12(1):109–138, 2001.
- [Nes05] Yu. Nesterov. Smooth minimization of non-smooth functions. *Math. Program.*, 103(1, Ser. A):127–152, 2005.
- [Nes16] Yu. Nesterov. Private communication, 2016. Les Houches, France.
- [NW06] J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer, New York, 2nd edition, 2006.
- [Ove88] M.L. Overton. On minimizing the maximum eigenvalue of a symmetric matrix. *SIAM J. Matrix Anal. Appl.*, 9:256–268, 1988.
- [Pow76] M. J. D. Powell. Some global convergence properties of a variable metric algorithm for minimization without exact line searches. In *Nonlinear Programming*, pages 53–72, Providence, 1976. Amer. Math. Soc. SIAM-AMS Proc., Vol. IX.
- [RFP10] Benjamin Recht, Maryam Fazel, and Pablo A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501, 2010.
- [Van19] Lieven Vandenberghe. Optimization methods for large-scale systems, 2019. <http://www.seas.ucla.edu/~vandenbe/236C/lectures/smoothing.pdf>, Lecture Notes for ECE236C.
- [XW17] Yuchen Xie and Andreas Waechter. On the convergence of BFGS on a class of piecewise linear non-smooth functions, December 2017. arXiv:1712.08571.