
A Quasi-Newton Approach to Nonsmooth Convex Optimization

Jin Yu

S.V. N. Vishwanathan

Simon Günter

Nicol N. Schraudolph

Canberra Research Laboratory, NICTA, Canberra, Australia

Research School of Information Sciences & Engineering, Australian National University, Canberra, Australia

JIN.YU@ANU.EDU.AU

SVN.VISHWANATHAN@NICTA.COM.AU

GUENTER_SIMON@HOTMAIL.COM

NIC@SCHRAUDOLPH.ORG

Abstract

We extend the well-known BFGS quasi-Newton method and its limited-memory variant LBFGS to the optimization of nonsmooth convex objectives. This is done in a rigorous fashion by generalizing three components of BFGS to subdifferentials: The local quadratic model, the identification of a descent direction, and the Wolfe line search conditions. We apply the resulting subLBFGS algorithm to L_2 -regularized risk minimization with binary hinge loss, and its direction-finding component to L_1 -regularized risk minimization with logistic loss. In both settings our generic algorithms perform comparable to or better than their counterparts in specialized state-of-the-art solvers.

1. Introduction

The (L)BFGS quasi-Newton method (Nocedal and Wright, 1999) is widely regarded as the workhorse of smooth nonlinear optimization due to its combination of computational efficiency with good asymptotic convergence. Given a smooth objective function $J : \mathbb{R}^d \rightarrow \mathbb{R}$ and a current iterate $\mathbf{w}_t \in \mathbb{R}^d$, BFGS forms a local quadratic model of J :

$$Q_t(\mathbf{p}) := J(\mathbf{w}_t) + \frac{1}{2} \mathbf{p}^\top \mathbf{B}_t^{-1} \mathbf{p} + \nabla J(\mathbf{w}_t)^\top \mathbf{p}, \quad (1)$$

where $\mathbf{B}_t \succ 0$ is a positive-definite estimate of the inverse Hessian of J . Minimizing $Q_t(\mathbf{p})$ gives the quasi-Newton direction

$$\mathbf{p}_t := -\mathbf{B}_t \nabla J(\mathbf{w}_t), \quad (2)$$

Appearing in *Proceedings of the 25th International Conference on Machine Learning*, Helsinki, Finland, 2008. Copyright 2008 by the author(s)/owner(s).

which is used for the parameter update:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \eta_t \mathbf{p}_t. \quad (3)$$

The step size $\eta_t \in \mathbb{R}_+$ is normally determined by a line search obeying the Wolfe conditions:

$$\begin{aligned} J(\mathbf{w}_{t+1}) &\leq J(\mathbf{w}_t) + c_1 \eta_t \nabla J(\mathbf{w}_t)^\top \mathbf{p}_t \\ \text{and } \nabla J(\mathbf{w}_{t+1})^\top \mathbf{p}_t &\geq c_2 \nabla J(\mathbf{w}_t)^\top \mathbf{p}_t, \end{aligned} \quad (4)$$

with $0 < c_1 < c_2 < 1$. The matrix \mathbf{B}_t is then modified via the incremental rank-two update

$$\mathbf{B}_{t+1} = (\mathbf{I} - \varrho_t \mathbf{s}_t \mathbf{y}_t^\top) \mathbf{B}_t (\mathbf{I} - \varrho_t \mathbf{y}_t \mathbf{s}_t^\top) + \varrho_t \mathbf{s}_t \mathbf{s}_t^\top, \quad (5)$$

where $\mathbf{s}_t := \mathbf{w}_{t+1} - \mathbf{w}_t$ and $\mathbf{y}_t := \nabla J(\mathbf{w}_{t+1}) - \nabla J(\mathbf{w}_t)$ denote the most recent step along the optimization trajectory in parameter and gradient space, respectively, and $\varrho_t := (\mathbf{y}_t^\top \mathbf{s}_t)^{-1}$. Given a descent direction \mathbf{p}_t , the Wolfe conditions ensure that $(\forall t) \mathbf{s}_t^\top \mathbf{y}_t > 0$ and hence $\mathbf{B}_0 \succ 0 \implies (\forall t) \mathbf{B}_t \succ 0$.

Limited-memory BFGS (LBFGS) is a variant of BFGS designed for solving large-scale optimization problems where the $O(d^2)$ cost of storing and updating \mathbf{B}_t would be prohibitively expensive. LBFGS approximates the quasi-Newton direction directly from the last m pairs of \mathbf{s}_t and \mathbf{y}_t via a matrix-free approach. This reduces the cost to $O(md)$ space and time per iteration, with m freely chosen (Nocedal and Wright, 1999).

Smoothness of the objective function is essential for standard (L)BFGS because both the local quadratic model (1) and the Wolfe conditions (4) require the existence of the gradient ∇J at every point. Even though nonsmooth convex functions are differentiable everywhere except on a set of Lebesgue measure zero (Hiriart-Urruty and Lemaréchal, 1993), in practice (L)BFGS often fails to converge on such problems (Lukšan and Vlček, 1999; Haarala, 2004). Various subgradient-based approaches, such as subgradient descent (Nedich and Bertsekas, 2000) or bundle methods (Teo et al., 2007), are therefore preferred.

Although a convex function might not be differentiable everywhere, a subgradient always exists. Let \mathbf{w} be a point where a convex function J is finite. Then a subgradient is the normal vector of any tangential supporting hyperplane of J at \mathbf{w} . Formally, \mathbf{g} is called a subgradient of J at \mathbf{w} if and only if

$$J(\mathbf{w}') \geq J(\mathbf{w}) + (\mathbf{w}' - \mathbf{w})^\top \mathbf{g} \quad \forall \mathbf{w}'. \quad (6)$$

The set of all subgradients at a point is called the subdifferential, and is denoted by $\partial J(\mathbf{w})$. If this set is not empty then J is said to be *subdifferentiable at \mathbf{w}* . If it contains exactly one element, *i.e.*, $\partial J(\mathbf{w}) = \{\nabla J(\mathbf{w})\}$, then J is *differentiable at \mathbf{w}* .

In this paper we systematically modify the standard (L)BFGS algorithm so as to make it amenable to subgradients. This results in sub(L)BFGS, a new subgradient quasi-Newton method which is applicable to a wide variety of nonsmooth convex optimization problems encountered in machine learning.

In the next section we describe our new algorithm generically, before we discuss its application to L_2 -regularized risk minimization with hinge loss in Section 3. Section 4 compares and contrasts our work with other recent efforts in this area. Encouraging experimental results are reported in Section 5. We conclude with an outlook and discussion in Section 6.

2. Subgradient BFGS Method

We modify the standard BFGS algorithm to derive our new algorithm (subBFGS, Algorithm 1) for nonsmooth convex optimization. These modifications can be grouped into three areas, which we elaborate on in turn: generalizing the local quadratic model, finding a descent direction, and finding a step size that obeys a subgradient reformulation of the Wolfe conditions.

2.1. Generalizing the Local Quadratic Model

Recall that BFGS assumes the objective function J is differentiable everywhere, so that at the current iterate \mathbf{w}_t we can construct a local quadratic model (1) of $J(\mathbf{w}_t)$. For a nonsmooth objective function, such a model becomes ambiguous at non-differentiable points (Figure 1). To resolve the ambiguity, we could simply replace the gradient $\nabla J(\mathbf{w}_t)$ in (1) with some subgradient $\mathbf{g}_t \in \partial J(\mathbf{w}_t)$. However, as will be discussed later, the resulting quasi-Newton direction $\mathbf{p}_t := -\mathbf{B}_t \mathbf{g}_t$ is not necessarily a descent direction. To address this fundamental modeling problem, we first generalize the

Algorithm 1 SUBGRADIENT BFGS (SUBBFGS)

```

1: Initialize:  $t := 0, \mathbf{w}_0 = \mathbf{0}, \mathbf{B}_0 = \mathbf{I}$ ;
2: Set direction-finding stopping tolerances  $\epsilon, k_{\max} \in \mathbb{R}_+$ ;
3: Compute subgradient  $\mathbf{g}_0 \in \partial J(\mathbf{w}_0)$ ;
4: while not converged do
5:    $\mathbf{p}_t = \text{descentDirection}(\mathbf{g}_t, \epsilon, k_{\max})$ ; (Algorithm 2)
6:   if  $\mathbf{p}_t = \text{failure}$  then
7:     Return  $\mathbf{w}_t$ ;
8:   end if
9:   Find  $\eta_t$  that obeys (14); (e.g., Algorithm 3)
10:   $\mathbf{s}_t = \eta_t \mathbf{p}_t$ ;
11:   $\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{s}_t$ ;
12:  Compute subgradient  $\mathbf{g}_{t+1} \in \partial J(\mathbf{w}_{t+1})$ ;
13:   $\mathbf{y}_t = \mathbf{g}_{t+1} - \mathbf{g}_t$ ;
14:  Update  $\mathbf{B}_{t+1}$  via (5);
15:   $t := t + 1$ ;
16: end while

```

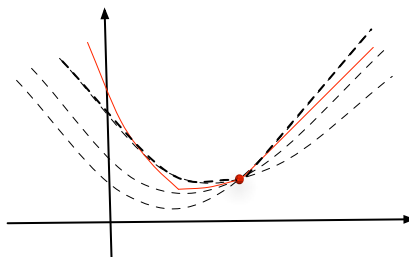


Figure 1. Quadratic models (dashed) vs. tightest pseudo-quadratic fit (7) (bold dashes) to the objective function (solid line) at a subdifferentiable point (solid disk).

local quadratic model as follows:

$$Q_t(\mathbf{p}) := J(\mathbf{w}_t) + M_t(\mathbf{p}), \quad \text{where} \\ M_t(\mathbf{p}) := \frac{1}{2} \mathbf{p}^\top \mathbf{B}_t^{-1} \mathbf{p} + \sup_{\mathbf{g} \in \partial J(\mathbf{w}_t)} \mathbf{g}^\top \mathbf{p}. \quad (7)$$

Note that where J is differentiable, (7) reduces to the familiar BFGS quadratic model (1). At non-differentiable points, however, the model is no longer quadratic, as the supremum may be attained at different elements of $\partial J(\mathbf{w}_t)$ for different directions \mathbf{p} . Instead it can be viewed as the tightest pseudo-quadratic fit to J at \mathbf{w}_t (Figure 1).

Ideally, we would like to minimize $Q_t(\mathbf{p})$, or equivalently $M_t(\mathbf{p})$, in (7) to obtain the best search direction,

$$\mathbf{p}^* := \underset{\mathbf{p} \in \mathbb{R}^d}{\text{arginf}} M_t(\mathbf{p}). \quad (8)$$

This is generally intractable due to the presence of a supremum over the entire subdifferential set $\partial J(\mathbf{w}_t)$. In many machine learning problems, however, the set $\partial J(\mathbf{w}_t)$ has some special structure that simplifies calculation of the supremum in (7). In what follows, we develop an iteration that is guaranteed to find a quasi-Newton descent direction, assuming an oracle that supplies $\text{argsup}_{\mathbf{g} \in \partial J(\mathbf{w}_t)} \mathbf{g}^\top \mathbf{p}$ for a given direction $\mathbf{p} \in \mathbb{R}^d$. In Section 3.1 we provide an efficient im-

Algorithm 2 $\mathbf{p}_t = \text{descentDirection}(\mathbf{g}^{(1)}, \epsilon, k_{\max})$

input subgradient $\mathbf{g}^{(1)} \in \partial J(\mathbf{w}_t)$,
tolerance $\epsilon \in \mathbb{R}_+$, iteration limit k_{\max} ;
output descent direction \mathbf{p}_t ;
1: Initialize: $i := 1$, $\bar{\mathbf{g}}^{(1)} = \mathbf{g}^{(1)}$, $\mathbf{p}^{(1)} = -\mathbf{B}_t \mathbf{g}^{(1)}$;
2: $\mathbf{g}^{(2)} = \text{argsup}_{\mathbf{g} \in \partial J(\mathbf{w}_t)} \mathbf{g}^\top \mathbf{p}^{(1)}$;
3: Calculate $\epsilon^{(1)}$ via (13);
4: **while** $(\mathbf{g}^{(i+1)})^\top \mathbf{p}^{(i)} > 0$ or $\epsilon^{(i)} > \epsilon$ and $i < k_{\max}$ **do**
5: $\mu^* := \min \left[1, \frac{(\mathbf{g}^{(i+1)} - \bar{\mathbf{g}}^{(i)})^\top \mathbf{p}^{(i)}}{(\mathbf{g}^{(i+1)} - \bar{\mathbf{g}}^{(i)})^\top \mathbf{B}_t (\mathbf{g}^{(i+1)} - \bar{\mathbf{g}}^{(i)})} \right]$;
6: $\bar{\mathbf{g}}^{(i+1)} = (1 - \mu^*) \bar{\mathbf{g}}^{(i)} + \mu^* \mathbf{g}^{(i+1)}$;
7: $\mathbf{p}^{(i+1)} = (1 - \mu^*) \mathbf{p}^{(i)} - \mu^* \mathbf{B}_t \mathbf{g}^{(i+1)}$;
8: $\mathbf{g}^{(i+2)} = \text{argsup}_{\mathbf{g} \in \partial J(\mathbf{w}_t)} \mathbf{g}^\top \mathbf{p}^{(i+1)}$;
9: Calculate $\epsilon^{(i+1)}$ via (13);
10: $i := i + 1$;
11: **end while**
12: **if** $(\mathbf{g}^{(i+1)})^\top \mathbf{p}^{(i)} > 0$ **then**
13: **return** failure;
14: **else**
15: **return** $\text{argmin}_{j \leq i} M_t(\mathbf{p}^{(j)})$.
16: **end if**

plementation of such an oracle for L_2 -regularized risk minimization with the hinge loss.

2.2. Finding a Descent Direction

A direction \mathbf{p}_t is a descent direction if and only if $\mathbf{g}^\top \mathbf{p}_t < 0 \forall \mathbf{g} \in \partial J(\mathbf{w}_t)$ (Belloni, 2005), or equivalently

$$\sup_{\mathbf{g} \in \partial J(\mathbf{w}_t)} \mathbf{g}^\top \mathbf{p}_t < 0. \quad (9)$$

In particular, for a smooth convex function the quasi-Newton direction (2) is always a descent direction because $\nabla J(\mathbf{w}_t)^\top \mathbf{p}_t = -\nabla J(\mathbf{w}_t)^\top \mathbf{B}_t \nabla J(\mathbf{w}_t) < 0$ holds due to the positivity of \mathbf{B}_t .

For nonsmooth functions, however, the quasi-Newton direction $\mathbf{p}_t := -\mathbf{B}_t \mathbf{g}_t$ for a given $\mathbf{g}_t \in \partial J(\mathbf{w}_t)$ may not fulfill the descent condition (9), making it impossible to find a step that obeys (4), thus causing a failure of the line search. We now present an iterative approach to finding a quasi-Newton *descent* direction.

Inspired by bundle methods (Teo et al., 2007), we build the following convex lower bound on $M_t(\mathbf{p})$:

$$M_t^{(i)}(\mathbf{p}) := \frac{1}{2} \mathbf{p}^\top \mathbf{B}_t^{-1} \mathbf{p} + \sup_{j \leq i} \mathbf{g}^{(j)\top} \mathbf{p}, \quad (10)$$

where $i, j \in \mathbb{N}$. Given a $\mathbf{p}^{(i)} \in \mathbb{R}^d$ the lower bound (10) is successively tightened by computing

$$\mathbf{g}^{(i+1)} := \text{argsup}_{\mathbf{g} \in \partial J(\mathbf{w}_t)} \mathbf{g}^\top \mathbf{p}^{(i)}, \quad (11)$$

such that $M_t^{(i)}(\mathbf{p}) \leq M_t^{(i+1)}(\mathbf{p}) \leq M_t(\mathbf{p}) \forall \mathbf{p} \in \mathbb{R}^d$. Here we set $\mathbf{g}^{(1)} \in \partial J(\mathbf{w}_t)$, and assume that $\mathbf{g}^{(i+1)}$ is

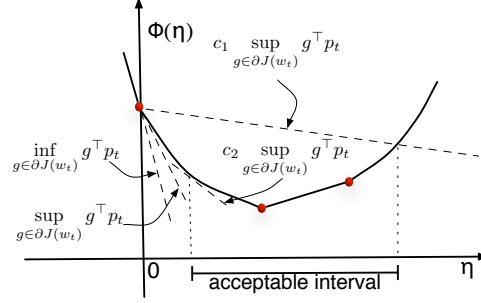


Figure 2. Geometric interpretation of the subgradient Wolfe conditions (14). Solid disks are subdifferentiable points; the slopes of dashed lines are indicated.

provided by an oracle. To solve $\inf_{\mathbf{p} \in \mathbb{R}^d} M_t^{(i)}(\mathbf{p})$, we rewrite it as a constrained optimization problem:

$$\inf_{\mathbf{p}, \xi} \left(\frac{1}{2} \mathbf{p}^\top \mathbf{B}_t^{-1} \mathbf{p} + \xi \right) \text{ s.t. } \mathbf{g}^{(j)\top} \mathbf{p} \leq \xi \quad \forall j \leq i. \quad (12)$$

This problem can be solved exactly via quadratic programming, but doing so may incur substantial computational expense. Instead we adopt an alternative approach (Algorithm 2) which does not solve $\inf_{\mathbf{p} \in \mathbb{R}^d} M_t^{(i)}(\mathbf{p})$ to optimality. The key idea is to write the proposed descent direction at iteration $i + 1$ as a convex combination of $\mathbf{p}^{(i)}$ and $-\mathbf{B}_t \mathbf{g}^{(i+1)}$. The optimal combination coefficient μ^* can be computed exactly (Step 5 of Algorithm 2) using an argument based on maximizing dual progress. Finally, to derive an implementable stopping criterion, we define $\epsilon^{(i)}$ to be

$$\min_{j \leq i} \left[\mathbf{p}^{(j)\top} \mathbf{g}^{(j+1)} - \frac{1}{2} (\mathbf{p}^{(j)\top} \bar{\mathbf{g}}^{(j)} + \mathbf{p}^{(i)\top} \bar{\mathbf{g}}^{(i)}) \right], \quad (13)$$

where $\bar{\mathbf{g}}^{(i)}$ is an aggregated subgradient (Step 6 of Algorithm 2) which lies in the convex hull of $\mathbf{g}^{(j)} \in \partial J(\mathbf{w}_t) \forall j \leq i$. $\epsilon^{(i)}$ is monotonically decreasing, and upper bounds the distance from the optimal value of the dual of $M_t(\mathbf{p})$, leading us to a practical stopping criterion (Step 4 of Algorithm 2) for our direction-finding procedure. Yu et al. (2008) provide details, and prove that Algorithm 2 converges to the optimal dual objective value with precision ϵ at an $O(1/\epsilon)$ rate.

2.3. Subgradient Line Search

Given the current iterate \mathbf{w}_t and a search direction \mathbf{p}_t , the task of a line search is to find a step size $\eta \in \mathbb{R}_+$ which decreases the objective function along the line $\mathbf{w}_t + \eta \mathbf{p}_t$, i.e., $J(\mathbf{w}_t + \eta \mathbf{p}_t) =: \Phi(\eta)$. The Wolfe conditions (4) are used in line search routines to enforce a sufficient decrease in the objective value, and to exclude unnecessarily small step sizes (Nocedal and Wright, 1999). However, the original Wolfe conditions require the objective function to be smooth. To extend

them to nonsmooth convex problems, we propose the following subgradient reformulation:

$$J(\mathbf{w}_{t+1}) \leq J(\mathbf{w}_t) + c_1 \eta_t \sup_{\mathbf{g} \in \partial J(\mathbf{w}_t)} \mathbf{g}^\top \mathbf{p}_t$$

$$\text{and } \sup_{\mathbf{g}' \in \partial J(\mathbf{w}_{t+1})} \mathbf{g}'^\top \mathbf{p}_t \geq c_2 \sup_{\mathbf{g} \in \partial J(\mathbf{w}_t)} \mathbf{g}^\top \mathbf{p}_t, \quad (14)$$

where $0 < c_1 < c_2 < 1$. Figure 2 illustrates how these conditions enforce acceptance of non-trivial step sizes that decrease the objective value. Yu et al. (2008) formally show that for any given descent direction we can always find a positive step size that satisfies (14).

2.4. Limited-Memory Subgradient BFGS

It is straightforward to implement an LBFGS variant of our subBFGS algorithm: We simply modify Algorithms 1 and 2 to compute all products of \mathbf{B}_t with a vector by means of the standard LBFGS matrix-free scheme (Nocedal and Wright, 1999).

3. sub(L)BFGS Implementation for L_2 -Regularized Risk Minimization

Many machine learning algorithms can be viewed as minimizing the L_2 -regularized risk

$$J(\mathbf{w}) := \frac{c}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n l(\mathbf{w}^\top \mathbf{x}_i, z_i), \quad (15)$$

where $\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^d$ are the training instances, $z_i \in \mathcal{Z} \subseteq \mathbb{R}$ the corresponding labels, and the loss l is a non-negative convex function of \mathbf{w} which measures the discrepancy between z_i and the predictions arising from \mathbf{w} via $\mathbf{w}^\top \mathbf{x}_i$. A loss function commonly used for binary classification is the hinge loss

$$l(\mathbf{w}^\top \mathbf{x}, z) := \max(0, 1 - z \mathbf{w}^\top \mathbf{x}), \quad (16)$$

where $z \in \{\pm 1\}$. L_2 -regularized risk minimization with binary hinge loss is a convex but nonsmooth optimization problem; in this section we show how sub(L)BFGS (Algorithm 1) can be applied to it.

Differentiating (15) after plugging in (16) yields

$$\partial J(\mathbf{w}) = c \mathbf{w} - \frac{1}{n} \sum_{i=1}^n \beta_i z_i \mathbf{x}_i = \bar{\mathbf{w}} - \frac{1}{n} \sum_{i \in \mathcal{M}} \beta_i z_i \mathbf{x}_i, \quad (17)$$

where $\bar{\mathbf{w}} := c \mathbf{w} - \frac{1}{n} \sum_{i \in \mathcal{E}} z_i \mathbf{x}_i$ and

$$\beta_i := \begin{cases} 1 & \text{if } i \in \mathcal{E}, \quad \mathcal{E} := \{i : 1 - z_i \mathbf{w}^\top \mathbf{x}_i > 0\}, \\ [0, 1] & \text{if } i \in \mathcal{M}, \quad \mathcal{M} := \{i : 1 - z_i \mathbf{w}^\top \mathbf{x}_i = 0\}, \\ 0 & \text{if } i \in \mathcal{W}, \quad \mathcal{W} := \{i : 1 - z_i \mathbf{w}^\top \mathbf{x}_i < 0\}. \end{cases}$$

\mathcal{E} , \mathcal{M} , and \mathcal{W} index the set of points which are in error, on the margin, and well-classified, respectively.

3.1. Realizing the Direction-Finding Method

Recall that our sub(L)BFGS algorithm requires an oracle that provides $\operatorname{argsup}_{\mathbf{g} \in \partial J(\mathbf{w}_t)} \mathbf{g}^\top \mathbf{p}$ for a given direction \mathbf{p} . For L_2 -regularized risk minimization with binary hinge loss we can implement such an oracle at computational cost linear in the number $|\mathcal{M}_t|$ of current marginal points. (Normally $|\mathcal{M}_t| \ll n$.) Towards this end we use (17) to obtain

$$\begin{aligned} \sup_{\mathbf{g} \in \partial J(\mathbf{w}_t)} \mathbf{g}^\top \mathbf{p} &= \sup_{\beta_i, i \in \mathcal{M}_t} \left(\bar{\mathbf{w}}_t - \frac{1}{n} \sum_{i \in \mathcal{M}_t} \beta_i z_i \mathbf{x}_i \right)^\top \mathbf{p} \\ &= \bar{\mathbf{w}}_t^\top \mathbf{p} - \frac{1}{n} \sum_{i \in \mathcal{M}_t} \inf_{\beta_i \in [0, 1]} \beta_i z_i \mathbf{x}_i^\top \mathbf{p}. \end{aligned} \quad (18)$$

Since for a given \mathbf{p} the first term of the right-hand side of (18) is a constant, the supremum is attained when we set $\beta_i \forall i \in \mathcal{M}_t$ via the following strategy:

$$\beta_i := \begin{cases} 0 & \text{if } z_i \mathbf{x}_i^\top \mathbf{p}_t \geq 0, \\ 1 & \text{if } z_i \mathbf{x}_i^\top \mathbf{p}_t < 0. \end{cases} \quad (19)$$

3.2. Implementing the Line Search

The one-dimensional convex function Φ obtained by restricting (15) to a line can be evaluated efficiently. To see this, rewrite the objective as

$$J(\mathbf{w}) := \frac{c}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \mathbf{1}^\top \max(\mathbf{0}, \mathbf{1} - \mathbf{z} \cdot \mathbf{X} \mathbf{w}), \quad (20)$$

where $\mathbf{0}$ and $\mathbf{1}$ are column vectors of zeros and ones, respectively, \cdot denotes the Hadamard (component-wise) product, and $\mathbf{z} \in \mathbb{R}^n$ collects correct labels corresponding to each row of data in $\mathbf{X} := [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$. Given a search direction \mathbf{p}_t at an iterate \mathbf{w}_t , this allows us to write

$$\begin{aligned} \Phi(\eta) := J(\mathbf{w}_t + \eta \mathbf{p}_t) &= \frac{1}{n} \delta(\eta)^\top [\mathbf{1} - (\mathbf{f} + \eta \Delta \mathbf{f})] \\ &\quad + \frac{c}{2} \|\mathbf{w}_t\|^2 + c \eta \mathbf{w}_t^\top \mathbf{p}_t + \frac{c \eta^2}{2} \|\mathbf{p}_t\|^2 \end{aligned} \quad (21)$$

where $\mathbf{f} := \mathbf{z} \cdot \mathbf{X} \mathbf{w}_t$, $\Delta \mathbf{f} := \mathbf{z} \cdot \mathbf{X} \mathbf{p}_t$, and

$$\delta_i(\eta) := \begin{cases} 1 & \text{if } f_i + \eta \Delta f_i < 1, \\ [0, 1] & \text{if } f_i + \eta \Delta f_i = 1, \\ 0 & \text{if } f_i + \eta \Delta f_i > 1 \end{cases} \quad (22)$$

for $1 \leq i \leq n$. We cache \mathbf{f} and $\Delta \mathbf{f}$, expending $O(nd)$ computational effort. We also cache $\frac{c}{2} \|\mathbf{w}_t\|^2$, $c \mathbf{w}_t^\top \mathbf{p}_t$, and $\frac{c}{2} \|\mathbf{p}_t\|^2$, each of which requires $O(n)$ work. The evaluation of $\delta(\eta)$ and its inner product with $\mathbf{1} - (\mathbf{f} + \eta \Delta \mathbf{f})$ both take $O(n)$ effort. All other terms in (21) can be computed in constant time, thus reducing the amortized cost of evaluating $\Phi(\eta)$ to $O(n)$. We are now in a position to introduce an exact line search which takes advantage of this scheme.

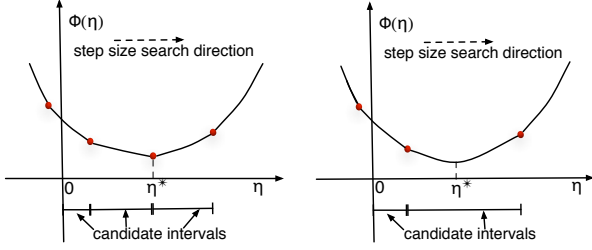


Figure 3. Nonsmooth convex function Φ of step size η . Solid disks are subdifferentiable points; the optimal η^* falls on such a point (left), or between two such points (right).

3.2.1. EXACT LINE SEARCH

Differentiating (21) with respect to η and setting the gradient to zero shows that $\eta^* := \operatorname{argmin}_{\eta} \Phi(\eta)$ satisfies $\eta^* = (\boldsymbol{\delta}(\eta^*)^\top \Delta \mathbf{f} / n - c \mathbf{w}_t^\top \mathbf{p}_t) / (c \|\mathbf{p}_t\|^2)$. It is easy to verify that $\Phi(\eta)$ is piecewise quadratic, and differentiable everywhere except at $\eta_i := (1 - f_i) / \Delta f_i$, where it becomes subdifferentiable. At these points an element of the indicator function $\boldsymbol{\delta}(\eta)$ (22) changes from 0 to 1 or vice versa; otherwise $\boldsymbol{\delta}(\eta)$ remains constant. Thus for a smooth interval (η_a, η_b) between subdifferentiable points η_a and η_b (cf. Figure 3), if the candidate step

$$\eta_{a,b}^* = \frac{\boldsymbol{\delta}(\eta')^\top \Delta \mathbf{f} / n - c \mathbf{w}_t^\top \mathbf{p}_t}{c \|\mathbf{p}_t\|^2}, \quad \eta' \in (\eta_a, \eta_b) \quad (23)$$

lies in the interval, it is optimal: $\eta_{a,b}^* \in [\eta_a, \eta_b] \Rightarrow \eta^* = \eta_{a,b}^*$. Otherwise, the interval boundary is optimal if its subdifferential contains zero: $0 \in \partial \Phi(\eta_a) \Rightarrow \eta^* = \eta_a$. Sorting the subdifferentiable points η_i facilitates efficient search over intervals; see Algorithm 3 for details.

4. Related Work

Lukšan and Vlček (1999) propose an extension of BFGS to nonsmooth convex problems. Their algorithm samples gradients around non-differentiable points in order to obtain a descent direction. In many machine learning problems evaluating the objective function and its gradient is very expensive. Therefore, our direction-finding algorithm (Algorithm 2) repeatedly samples subgradients from the set $\partial J(\mathbf{w})$ via the oracle, which is computationally more efficient.

Recently, Andrew and Gao (2007) introduced a variant of nonsmooth BFGS, the Orthant-Wise Limited-memory Quasi-Newton (OWL-QN) algorithm, suitable for optimizing L_1 -regularized log-linear models:

$$J(\mathbf{w}) := c \|\mathbf{w}\|_1 + \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-z_i \mathbf{w}^\top \mathbf{x}_i}), \quad (24)$$

where the logistic loss is smooth, but the regularizer

Algorithm 3 $\eta = \text{linesearch}(\mathbf{w}_t, \mathbf{p}_t, c, \mathbf{f}, \Delta \mathbf{f})$

input $\mathbf{w}_t, \mathbf{p}_t, c, \mathbf{f}$, and $\Delta \mathbf{f}$ as in (21);
output step size η ;
1: $b = c \mathbf{w}_t^\top \mathbf{p}_t$, $h = c \|\mathbf{p}_t\|^2$;
2: $n = \text{length}(\mathbf{f})$, $j := 1$;
3: $\boldsymbol{\alpha} := [(1 - \mathbf{f}) / \Delta \mathbf{f}, 0]$; (subdifferentiable points)
4: $\boldsymbol{\pi} = \text{sort}(\boldsymbol{\alpha})$; (index vector)
5: **while** $\alpha_{\pi_j} \leq 0$ **do**
6: $j := j + 1$;
7: **end while**
8: $\eta := \alpha_{\pi_j} / 2$;
9: **for** $i := 1$ to n **do**
10: $\delta_i := \begin{cases} 1 & \text{if } f_i + \eta \Delta f_i < 1; \\ 0 & \text{otherwise;} \end{cases}$
11: **end for**
12: $\varrho := \boldsymbol{\delta}^\top \Delta \mathbf{f} / n$;
13: **while** $j \leq \text{length}(\boldsymbol{\pi})$ **do**
14: $\eta := (\varrho - b) / h$; (candidate step)
15: **if** $\eta \in [\alpha_{\pi_{j-1}}, \alpha_{\pi_j}]$ **then**
16: **return** η ;
17: **else if** $\eta < \alpha_{\pi_{j-1}}$ **then**
18: **return** $\eta := \alpha_{\pi_{j-1}}$;
19: **else**
20: **repeat**
21: $\varrho := \begin{cases} \varrho - \Delta f_{\pi_j} / n & \text{if } \delta_{\pi_j} = 1, \\ \varrho + \Delta f_{\pi_j} / n & \text{otherwise;} \end{cases}$
22: $j := j + 1$;
23: **until** $\alpha_{\pi_j} \neq \alpha_{\pi_{j-1}}$
24: **end if**
25: **end while**

is only subdifferentiable at points where \mathbf{w} has zero elements. From the optimization viewpoint this objective is very similar to the L_2 -regularized hinge loss; the direction finding and line search methods that we discussed in Sections 3.1 and 3.2, respectively, can be applied to this problem with slight modifications.

OWL-QN is based on the observation that the L_1 regularizer is linear within any given orthant. Therefore, it maintains an approximation \mathbf{B}^{ow} to the inverse Hessian of the logistic loss, and uses an efficient scheme to select orthants for optimization. In fact, its success greatly depends on its direction-finding subroutine, which demands a specially chosen subgradient \mathbf{g}^{ow} (Andrew and Gao, 2007, Equation 4) to produce the quasi-Newton direction, $\mathbf{p}^{\text{ow}} = \pi(\mathbf{p}, \mathbf{g}^{\text{ow}})$, where $\mathbf{p} := -\mathbf{B}^{\text{ow}} \mathbf{g}^{\text{ow}}$ and the projection π returns a search direction by setting the i^{th} element of \mathbf{p} to zero whenever $p_i g_i^{\text{ow}} > 0$. As shown in Section 5, the direction-finding subroutine of OWL-QN can be replaced by Algorithm 2, which in turn makes the algorithm more robust to the choice of subgradients.

Many optimization techniques use past gradients to build a model of the objective function. Bundle method solvers like SVMStruct (Joachims, 2006) and BMRM (Teo et al., 2007) use them to lower-bound the

Table 1. Datasets, regularization constants c , direction-finding convergence criterion ϵ , and the overall number k of direction-finding iterations for L_1 -regularized logistic loss and L_2 -regularized hinge loss minimization tasks, respectively.

Dataset	Tr./Test Data	Dimen.	Density	c_{L_1}	ϵ_{L_1}	k_{L_1}	$k_{L_1\text{rand}}$	c_{L_2}	ϵ_{L_2}	k_{L_2}
Coverttype	522911/58101	54	22.22%	10^{-6}	10^{-5}	0	0	10^{-6}	10^{-8}	44
CCAT	781265/23149	47236	0.16%	10^{-6}	10^{-5}	356	467	10^{-4}	10^{-8}	66
Astro	29882/32487	99757	0.077%	10^{-5}	10^{-3}	1668	2840	$5 \cdot 10^{-5}$	10^{-8}	17
MNIST	60000/10000	780	19.22%	10^{-4}	10^{-5}	60	102	$1.4286 \cdot 10^{-6}$	10^{-8}	244

objective by a piecewise linear function which is minimized to obtain the next iterate. This fundamentally differs from the BFGS approach of using past gradients to approximate the (inverse) Hessian, hence building a quadratic model of the objective function.

Vojtěch and Sonnenburg (2007) speed up the convergence of a bundle method solver for the L_2 -regularized binary hinge loss. Their main idea is to perform a line search along the line connecting two successive iterates of a bundle method solver. Although developed independently, their line search algorithm is very reminiscent of the method we describe in Section 3.2.1.

5. Experiments

We now evaluate the performance of our subLBFGS algorithm, and compare it to other state-of-the-art nonsmooth optimization methods on L_2 -regularized hinge loss minimization. We also compare a variant of OWL-QN that uses our direction-finding routine to the original on L_1 -regularized logistic loss minimization.

Our experiments used four datasets: the Coverttype dataset of Blackard, Jock & Dean, CCAT from the Reuters RCV1 collection, the Astro-physics dataset of abstracts of scientific papers from the Physics ArXiv (Joachims, 2006), and the MNIST dataset of handwritten digits with two classes: even and odd digits. We used subLBFGS with a buffer of size $m = 15$ throughout. Table 1 summarizes our parameter settings, and reports the overall number of direction-finding iterations for all experiments. We followed the choices of Vojtěch and Sonnenburg (2007) for the L_2 regularization constants; for L_1 they were chosen from the set $10^{\{-6, -5, \dots, -1\}}$ to achieve the lowest test error.

On convex problems such as these every convergent optimizer will reach the same solution; comparing generalisation performance is therefore pointless. We combined training and test datasets to evaluate the convergence of each algorithm in terms of the objective function value *vs.* CPU seconds. All experiments were carried out on a Linux machine with dual 2.8 GHz Xeon processors with 4GB RAM.

5.1. L_2 -Regularized Hinge Loss

For our first set of experiments, we applied subLBFGS together with our exact line search (Algorithm 3) to the task of L_2 -regularized hinge loss minimization. Our control methods are the bundle method solver BMRM (Teo et al., 2007) and an optimized cutting plane algorithm, OCAS version 0.6.0 (Vojtěch and Sonnenburg, 2007), both of which demonstrated strong results on the L_2 -regularized hinge loss minimization in their corresponding papers.

Figure 4 shows that subLBFGS (solid) reaches the neighbourhood of the optimum (less than 10^{-3} away) noticeably (up to 7 times) faster than BMRM (dashed). As BMRM’s approximation to the objective function improves over the course of optimization, it gradually catches up with subLBFGS, though ultimately subLBFGS still converges faster on 3 out of 4 datasets. The performance of subLBFGS and OCAS (dash-dotted) are very similar: OCAS converges slightly faster than subLBFGS on the Astrophysics dataset but is outperformed by subLBFGS on the MNIST dataset.

5.2. L_1 -Regularized Logistic Loss

To demonstrate the utility of our direction-finding routine (Algorithm 2) in its own right, we plugged it into the OWL-QN algorithm (Andrew and Gao, 2007) as an alternative direction-finding method, such that $\mathbf{p}^{\text{ow}} = \text{descentDirection}(\mathbf{g}^{\text{ow}}, \epsilon, k_{\text{max}})$,¹ and compared this variant (denoted by OWL-QN*) with the original on L_1 -regularized logistic loss minimization.

Using the stopping criterion suggested by Andrew and Gao (2007), we run experiments until the averaged relative change in the objective value over the previous 5 iterations falls below 10^{-5} . Figure 5 shows only minor differences in convergence between the two algorithms.

To examine the algorithms’ sensitivity to the choice of subgradients, we also ran them with random subgradients (as opposed to the specially chosen \mathbf{g}^{ow} used be-

¹Note for the objective (24) it is trivial to construct an oracle that supplies $\text{argsup}_{\mathbf{g} \in \partial J(\mathbf{w}_t)} \mathbf{g}^\top \mathbf{p}$.

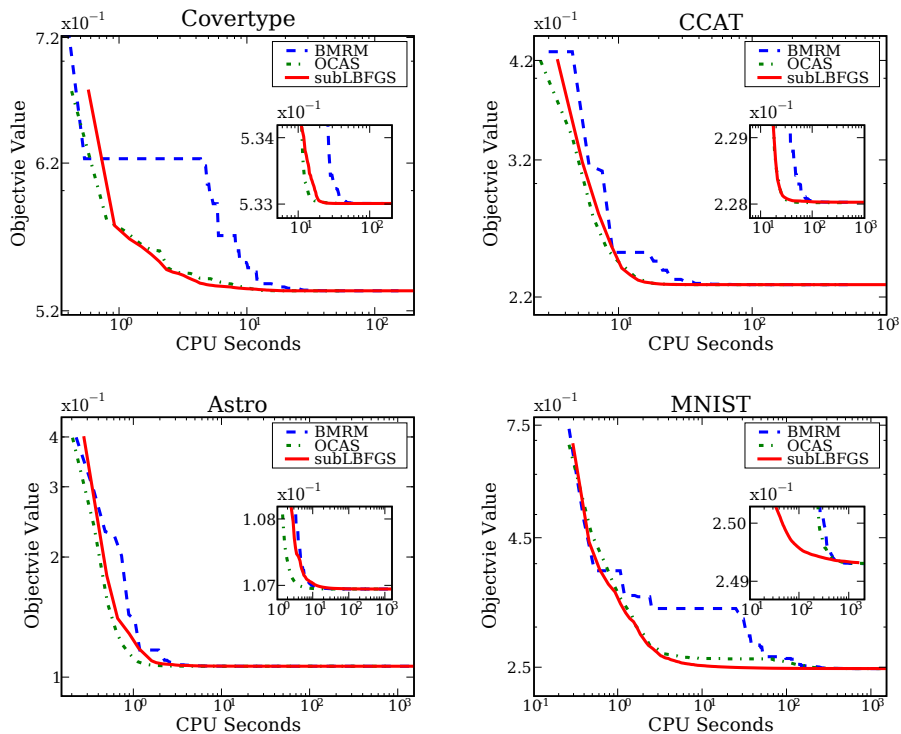


Figure 4. Objective function value *vs.* CPU seconds on L_2 -regularized hinge loss minimization tasks.

fore) fed to their corresponding direction-finding routines. OWL-QN relies heavily on its particular choice of subgradients, hence breaks down completely under these conditions: The only dataset where we could even plot its (poor) performance was Coverttype (dotted OWL-QN(2) line in Figure 5). Our direction-finding routine, by contrast, is self-correcting and thus not affected by this manipulation: The curves for OWL-QN*(2) (plotted for Coverttype in Figure 5) lie virtually on top of those for OWL-QN*. Table 1 shows that in this case more direction-finding iterations are needed, *i.e.*, $k_{L_1\text{rand}} \geq k_{L_1}$. This empirically confirms that as long as $\text{argsup}_{g \in \partial J(w_t)} g^\top p$ is given, Algorithm 2 can indeed be used as a canned quasi-Newton direction-finding routine.

6. Outlook and Discussion

We proposed an extension of BFGS suitable for handling nonsmooth problems often encountered in the machine learning context. As our experiments show, our algorithms are versatile and applicable to many problems, while their performance is comparable to if not better than that of their counterparts in custom-built solvers.

In some experiments we observe that subLBFGS initially makes rapid progress towards the solution but

slows down closer to the optimum. We hypothesize that initially its quadratic model allows subLBFGS to make rapid progress, but closer to the optimum it is no longer an accurate model of an objective function dominated by the nonsmooth hinges. We are therefore contemplating hybrid solvers which seamlessly switch between sub(L)BFGS and bundle solvers.

In this paper we applied subLBFGS to L_2 -regularized risk minimization with binary hinge loss. It can also be extended to deal with generalizations of the hinge loss, such as multi-class, multi-category, and ordinal regression problems; this is part of our ongoing research.

Finally, to put our contributions in perspective, recall that we modified three aspects of the standard BFGS algorithm, namely the quadratic model (Section 2.1), the descent direction finding (Section 2.2), and the line search (Section 2.3). Each of these modifications is versatile enough to be used as a component in other nonsmooth optimization algorithms. This not only offers the promise of improving existing algorithms, but may also help clarify connections between them. We hope that this will focus attention on those core sub-routines that need to be made more efficient in order to handle larger and larger datasets.

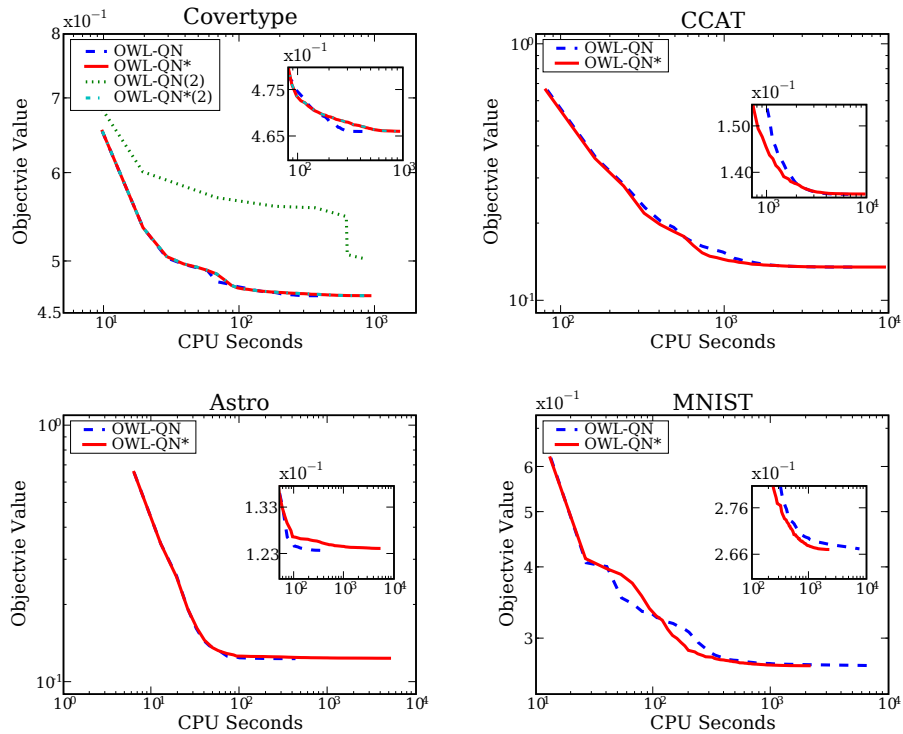


Figure 5. Objective function value vs. CPU seconds on L_1 -regularized logistic loss minimization tasks.

Acknowledgement

NICTA is funded by the Australian Government’s Backing Australia’s Ability and the Centre of Excellence programs. This work is also supported by the IST Program of the European Community, under the FP7 Network of Excellence, ICT-216886-NOE.

References

- G. Andrew and J. Gao. Scalable training of l_1 -regularized log-linear models. In *Proc. Intl. Conf. Machine Learning*, pages 33–40, New York, NY, USA, 2007. ACM.
- A. Belloni. Introduction to bundle methods. Technical report, Operation Research Center, M.I.T., 2005.
- M. Haarala. *Large-Scale Nonsmooth Optimization*. PhD thesis, University of Jyväskylä, 2004.
- J. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms, I and II*, volume 305 and 306. Springer-Verlag, 1993.
- T. Joachims. Training linear SVMs in linear time. In *Proc. ACM Conf. Knowledge Discovery and Data Mining (KDD)*. ACM, 2006.
- L. Lukšan and J. Vlček. Globally convergent variable metric method for convex nonsmooth unconstrained minimization. *Journal of Optimization Theory and Applications*, 102(3):593–613, 1999.
- A. Nedich and D. P. Bertsekas. Convergence rate of incremental subgradient algorithms. In S. Uryasev and P. M. Pardalos, editors, *Stochastic Optimization: Algorithms and Applications*, pages 263–304. Kluwer Academic Publishers, 2000.
- J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer Series in Operations Research. Springer, 1999.
- C. Teo, Q. Le, A. Smola, and S. Vishwanathan. A scalable modular convex solver for regularized risk minimization. In *Proc. ACM Conf. Knowledge Discovery and Data Mining (KDD)*. ACM, 2007.
- F. Vojtěch and S. Sonnenburg. Optimized cutting plane algorithm for support vector machines. Technical Report 1, Fraunhofer Institute FIRST, December 2007. <http://publica.fraunhofer.de/documents/N-66225.html>.
- J. Yu, S. V. N. Vishwanathan, S. Günter, and N. N. Schraudolph. A quasi-Newton approach to nonsmooth convex optimization. Technical Report arXiv:0804.3835, April 2008. <http://arxiv.org/pdf/0804.3835>.