

AN INEXACT VARIABLE METRIC PROXIMAL POINT ALGORITHM FOR GENERIC QUASI-NEWTON ACCELERATION*

HONGZHOU LIN[†], JULIEN MAIRAL[‡], AND ZAID HARCHAOU[§]

Abstract. We propose an inexact variable-metric proximal point algorithm to accelerate gradient-based optimization algorithms. The proposed scheme, called QNing, can notably be applied to incremental first-order methods such as the stochastic variance-reduced gradient descent algorithm and other randomized incremental optimization algorithms. QNing is also compatible with composite objectives, meaning that it has the ability to provide exactly sparse solutions when the objective involves a sparsity-inducing regularization. When combined with limited-memory BFGS rules, QNing is particularly effective at solving high-dimensional optimization problems while enjoying a worst-case linear convergence rate for strongly convex problems. We present experimental results where QNing gives significant improvements over competing methods for training machine learning methods on large samples and in high dimensions.

Key words. convex optimization, quasi-Newton, L-BFGS

AMS subject classifications. 90C25, 90C53

DOI. 10.1137/17M1125157

1. Introduction. Convex composite optimization arises in many scientific fields, such as image and signal processing or machine learning. It consists of minimizing a real-valued function composed of two convex terms:

$$(1) \quad \min_{x \in \mathbb{R}^d} \{f(x) \triangleq f_0(x) + \psi(x)\},$$

where f_0 is smooth with Lipschitz continuous derivatives, and ψ is a regularization function that is not necessarily differentiable. A typical example from the signal and image processing literature is the ℓ_1 -norm $\psi(x) = \|x\|_1$, which encourages sparse solutions [19, 40]; composite minimization also encompasses constrained minimization when considering extended-valued indicator functions ψ that may take the value $+\infty$ outside of a convex set \mathcal{C} and 0 inside (see [28]). In general, algorithms that are dedicated to composite optimization only require the ability to efficiently compute the proximal operator of ψ :

$$p_\psi(y) \triangleq \arg \min_{x \in \mathbb{R}^d} \left\{ \psi(x) + \frac{1}{2} \|x - y\|^2 \right\},$$

where $\|\cdot\|$ denotes the Euclidean norm. Note that when ψ is an indicator function the proximal operator corresponds to the simple Euclidean projection.

*Received by the editors April 10, 2017; accepted for publication (in revised form) January 22, 2019; published electronically May 28, 2019.

<http://www.siam.org/journals/siopt/29-2/M112515.html>

Funding: This work was supported by ERC grant SOLARIS (714381), a grant from the ANR (MACARON project ANR-14-CE23-0003-01), and the program “Learning in Machines and Brains” (CIFAR).

[†]Université Grenoble Alpes, INRIA, CNRS, Grenoble INP, LJK, Grenoble 38330, France. Current address: Computer Science and Artificial Intelligence Laboratory, MIT, 32 Vassar Street, Cambridge, MA 02139 (hongzhou@mit.edu).

[‡]Université Grenoble Alpes, INRIA, CNRS, Grenoble INP, LJK, Grenoble 38330, France (julien.mairal@inria.fr).

[§]Department of Statistics, University of Washington, Seattle, WA 98195-4322 (zaid@uw.edu).

To solve (1), significant efforts have been devoted to the following:

- (i) extending techniques for smooth optimization to deal with composite terms;
- (ii) exploiting the underlying structure of the problem, i.e.,
 - is f a finite sum of independent terms?
 - is ψ separable in different blocks of coordinates?
- (iii) exploiting the local curvature of the smooth term f to achieve faster convergence than gradient-based approaches when dimension d is large.

Typically, the first point is well understood in the context of optimal first-order methods (see [2, 48]), and the third point is tackled with effective heuristics such as the limited-memory BFGS (L-BFGS) algorithm when the problem is smooth [35, 49]. Yet, addressing all these challenges at the same time, which is precisely the focus of this paper, is difficult.

In particular, a problem of interest that initially motivated our work is that of empirical risk minimization (ERM); the problem arises in machine learning and can be formulated as the minimization of a composite function $f : \mathbb{R}^d \rightarrow \mathbb{R}$:

$$(2) \quad \min_{x \in \mathbb{R}^d} \left\{ f(x) \triangleq \frac{1}{n} \sum_{i=1}^n f_i(x) + \psi(x) \right\},$$

where the functions f_i are convex and smooth with Lipschitz continuous derivatives, and ψ is a composite term, possibly nonsmooth. The function f_i measures the fit of some model parameters x to a specific data point indexed by i , and ψ is a regularization penalty to prevent overfitting. To exploit the sum structure of f , a large number of randomized incremental gradient-based techniques have been proposed, such as stochastic average gradient (SAG) [56], SAGA [15], stochastic dual coordinate ascent (SDCA) [58], stochastic variance-reduced gradient descent (SVRG) [60], Finito [16], or minimization by incremental surrogate optimization (MISO) [38]. These approaches access a single gradient $\nabla f_i(x)$ at every iteration instead of the full gradient $(1/n) \sum_{i=1}^n \nabla f_i(x)$ and achieve lower computational complexity in expectation than optimal first-order methods [2, 48] under a few assumptions. Yet, these methods are unable to exploit the curvature of the objective function; indeed, this is also the case for variants that are accelerated in the sense of Nesterov [21, 33, 58].

To tackle (2), dedicated first-order methods are often the default choice in machine learning, but it is also known that standard quasi-Newton approaches can sometimes be surprisingly effective in the smooth case—that is, when $\psi = 0$ (see, e.g., [56] for extensive benchmarks). Since the dimension, d , of the problem is typically very large ($d \geq 10\,000$), “limited-memory” variants of these algorithms, such as L-BFGS, are necessary to achieve the desired scalability [35, 49]. The theoretical guarantees offered by L-BFGS are somewhat limited, meaning that it does not outperform accelerated first-order methods in terms of worst-case convergence rate and also that it is not guaranteed to correctly approximate the Hessian of the objective. Yet, L-BFGS remains one of the greatest practical successes of smooth optimization. Adapting L-BFGS to composite and structured problems, such as the finite sum of functions (2), has become increasingly important.

For instance, there have been several attempts to develop a proximal quasi-Newton method [10, 31, 54, 62]. These algorithms typically require the proximal operator of ψ to be computed many times with respect to a variable metric. Quasi-Newton steps have also been incorporated as local search steps into accelerated first-order methods to further enhance their numerical performance [24]. More related to our work, in [26] L-BFGS is combined with SVRG for minimizing smooth finite sums.

The scope of our approach is broader, beyond the case of SVRG. We present a generic quasi-Newton scheme, applicable to a large class of first-order methods for composite optimization including other incremental algorithms [15, 16, 38, 56, 58] and block coordinate descent methods [51, 52].

More precisely, the main contribution of this paper is a generic meta-algorithm, called QNing (the letters “Q” and “N” stand for quasi-Newton), which uses a given optimization method to solve a sequence of auxiliary problems up to some appropriate accuracy, resulting in faster global convergence in practice. QNing falls into the class of inexact proximal point algorithms with variable metric and *may be seen as applying a quasi-Newton algorithm with inexact (but accurate enough) gradients to the Moreau envelope of the objective*. As a result, (i) our approach is generic, as stated previously; (ii) despite the smoothing of the objective, the subproblems that we solve are composite ones, which may lead to exactly sparse iterates when a sparsity-inducing regularization, e.g., the ℓ_1 -norm, is involved; (iii) when used with L-BFGS rules, it admits a worst-case linear convergence rate for strongly convex problems similar to that of gradient descent (GD), which is typically the best guarantee obtained for L-BFGS schemes in the literature.

The idea of combining second-order or quasi-Newton methods with the Moreau envelope is in fact relatively old. It may be traced back to variable metric proximal bundle methods [14, 23, 41], which aim to incorporate curvature information into the bundle methods. Our approach revisits this principle with a *limited-memory variant* (to deal with large dimension d), a *simple line-search scheme*, *several warm-start strategies for the subproblems*, and a *global complexity analysis*, which is more relevant than convergence rates of the iterates, as the latter do not take into account the cost per iteration.

To demonstrate the effectiveness of our scheme in practice, we evaluate QNing on regularized logistic regression and regularized least squares, with smooth and non-smooth regularization penalties such as the elastic net [63]. We use large-scale machine learning data sets and show that QNing performs at least as well as the recently proposed accelerated incremental algorithm Catalyst [33] and other quasi-Newton baselines, such as proximal quasi-Newton methods [31] and stochastic L-BFGS [44], in all numerical experiments, and significantly outperforms them in many cases.

The paper is organized as follows: section 2 presents related work on quasi-Newton methods such as L-BFGS; we introduce QNing in section 3 and its convergence analysis in section 4; section 5 is devoted to numerical experiments; and section 6 concludes the paper.

2. Related work and preliminaries. The history of quasi-Newton methods can be traced back to the 1950s [6, 29, 50]. In practice, quasi-Newton methods often lead to significantly faster convergence than simpler gradient-based methods for solving smooth optimization problems [55]. Yet, a theoretical analysis of quasi-Newton methods that explains their impressive empirical behavior is still an open problem. We briefly review the well-known BFGS algorithm in section 2.1, and review its limited-memory variant [49] and a few recent extensions in section 2.2. Then, in section 2.3 we discuss earlier works that combine proximal point algorithm and quasi-Newton methods.

2.1. Quasi-Newton methods for smooth optimization. The most popular quasi-Newton methods are probably BFGS, named after its inventors (Broyden, Fletcher, Goldfarb, and Shanno), and its limited variant, L-BFGS [50]. These approaches will be the workhorses of the QNing meta-algorithm in practice. Consider

a smooth convex objective f to be minimized. The BFGS method constructs at iteration k a couple (x_k, B_k) with the following update:

$$(3) \quad x_{k+1} = x_k - \alpha_k B_k^{-1} \nabla f(x_k) \quad \text{and} \quad B_{k+1} = B_k - \frac{B_k s_k s_k^\top B_k}{s_k^\top B_k s_k} + \frac{y_k y_k^\top}{y_k^\top s_k},$$

where α_k is a suitable step size and

$$s_k = x_{k+1} - x_k, \quad y_k = \nabla f(x_{k+1}) - \nabla f(x_k).$$

The matrix B_k aims to approximate the Hessian matrix at iterate x_k . When f is strongly convex, the positive definiteness of B_k is guaranteed, as well as the condition $y_k^\top s_k > 0$, which ensures that (3) is well defined. The step size α_k is usually determined by a line-search strategy. For instance, applying Wolfe's line-search strategy provides a linear convergence rate for strongly convex objectives. Moreover, under the stronger conditions that the objective f is twice differentiable and its Hessian is Lipschitz continuous, the algorithm can asymptotically achieve a superlinear convergence rate [50].

However, when the dimension d is large, storing the d -by- d matrix B_k is infeasible. The limited-memory variant L-BFGS [49] overcomes this issue by restricting the matrix B_k to be low rank. More precisely, instead of storing the full matrix, a "generating list" of at most l pairs of vectors $\{(s_i^k, y_i^k)\}_{i=0, \dots, j}$ is kept in the memory. The low rank matrix B_k can then be recovered by performing the matrix update recursion in (3) involving all pairs of the generating list. Between iterations k and $k+1$, the generating list is incrementally updated by removing the oldest pair in the list (when $j = l$) and adding a new one. What makes the approach appealing is the ability to compute the matrix-vector product $H_k z = B_k^{-1} z$ with only $O(ld)$ floating-point operations for any vector z . This procedure relies entirely on a vector-vector product that does not explicitly construct the d -by- d matrix B_k or H_k . The price to pay is that superlinear convergence becomes out of reach.

L-BFGS is thus appropriate for high-dimensional problems (when d is large), but still requires the full gradient to be computed at each iteration, which may be cumbersome in the large sum setting (2). This motivated a stochastic counterpart of the quasi-Newton method (called stochastic quasi-Newton, or SQN) [9, 42, 57]. Unfortunately, directly substituting the full gradient $\nabla f(x_k)$ by its stochastic counterpart does not lead to a convergent scheme. Instead, the SQN method [9] uses updates with subsampled Hessian-vector products, which leads to a sublinear convergence rate. Later, a linearly convergent SQN algorithm was proposed by exploiting a variance reduction scheme [26, 44]. However, it is unclear how to extend these techniques to the composite setting.

2.2. Quasi-Newton methods for composite optimization. Different approaches have been proposed to extend quasi-Newton methods to composite optimization problems. A first approach consists in minimizing successive quadratic approximations, also called proximal quasi-Newton methods [10, 25, 30, 31, 36, 54]. More concretely, a local quadratic approximation q_k is minimized at each iteration:

$$(4) \quad q_k(x) \triangleq f_0(x_k) + \langle \nabla f_0(x_k), x - x_k \rangle + \frac{1}{2}(x - x_k)^\top B_k (x - x_k) + \psi(x),$$

where B_k is a Hessian approximation based on quasi-Newton methods. The minimizer of q_k provides a descent direction, which is subsequently used to build the next iterate.

However, a closed-form solution of (4) is usually not available since B_k changes over the iterations. Thus, one needs to apply an optimization algorithm to approximately solve (4). The composite structure of the subproblem naturally leads to the choice of a first-order optimization algorithm, such as a randomized coordinate descent algorithm. Then, superlinear complexity becomes out of reach since it requires the subproblems (4) to be solved with “high accuracy” [31]. The global convergence rate of this inexact variant was analyzed, for instance, in [54], where a sublinear convergence rate was obtained for convex problems; later, the analysis was extended to strongly convex problems in [36], where a linear convergence rate was achieved.

A second approach to extending quasi-Newton methods to composite optimization problems is based on smoothing techniques. More precisely, a quasi-Newton method is applied to a smoothed version of the objective. For instance, one may use the forward-backward envelope [4, 59]. The idea is to mimic forward-backward splitting methods and apply quasi-Newton steps instead of gradient steps on top of the envelope. Another well-known smoothing technique is to apply the Moreau–Yosida regularization [43, 61], which gives a smoothed function called the Moreau envelope. Then, applying quasi-Newton methods on it leads to the family of variable metric proximal point algorithms [7, 14, 22, 23]. Our method pursues this line of work by developing a practical inexact variant with global complexity guarantees.

2.3. Combining the proximal point algorithm and quasi-Newton methods. We briefly recall the definition of the Moreau envelope and its basic properties.

DEFINITION 1. *Given an objective function f and a smoothing parameter $\kappa > 0$, the Moreau envelope of f is the function F obtained by performing the infimal convolution*

$$(5) \quad F(x) \triangleq \min_{z \in \mathbb{R}^d} \left\{ f(z) + \frac{\kappa}{2} \|z - x\|^2 \right\}.$$

When f is convex, the subproblem defined in (5) is strongly convex, which provides a unique minimizer, called the *proximal point* of x , which we denote by $p(x)$.

PROPOSITION 1 (basic properties of the Moreau envelope). *If f is convex, the Moreau envelope F defined in (5) satisfies the following.*

1. F has the same minimum as f , i.e.,

$$\min_{x \in \mathbb{R}^d} F(x) = \min_{x \in \mathbb{R}^d} f(x),$$

and the solution sets of the above two problems coincide with each other.

2. F is continuously differentiable even when f is not, and

$$(6) \quad \nabla F(x) = \kappa(x - p(x)).$$

Moreover, the gradient ∇F is Lipschitz continuous with constant $L_F = \kappa$.

3. F is convex. Moreover, when f is μ -strongly convex with respect to the Euclidean norm, F is μ_F -strongly convex with $\mu_F = \frac{\mu\kappa}{\mu+\kappa}$.

4. F is upper-bounded by f . More precisely, for any $x \in \mathbb{R}^d$,

$$(7) \quad F(x) + \frac{1}{2\kappa} \|\nabla F(x)\|^2 \leq f(x).$$

Interestingly, F inherits all the convex properties of f and, more importantly, it is always continuously differentiable (see [32] for elementary proofs). Moreover, the

condition number of F is given by

$$(8) \quad q = \frac{L_F}{\mu_F} = \frac{\mu + \kappa}{\mu},$$

which may be adjusted by the regularization parameter κ . Then, a naive approach to overcome the nonsmoothness of the function f is to transfer the optimization problem to its Moreau envelope F . More concretely, we may apply an optimization algorithm to minimize F and use the obtained solution as a solution to the original problem, since both functions share the same minimizers. This yields the following well-known algorithm.

Proximal point algorithm. Consider the gradient descent method with constant step size $1/L_F = 1/\kappa$:

$$x_{k+1} = x_k - \frac{1}{\kappa} \nabla F(x_k).$$

By rewriting the gradient $\nabla F(x_k)$ as $\kappa(x_k - p(x_k))$, we obtain the proximal point algorithm [53]:

$$(9) \quad x_{k+1} = p(x_k) = \arg \min_{z \in \mathbb{R}^d} \left\{ f(z) + \frac{\kappa}{2} \|z - x_k\|^2 \right\}.$$

Accelerated proximal point algorithm. Since GD on F yields the proximal point algorithm, it is natural to apply an accelerated first-order method to get faster convergence. To that effect, Nesterov's algorithm [45] uses a two-stage update, along with a specific extrapolation parameter β_{k+1} ,

$$x_{k+1} = y_k - \frac{1}{\kappa} \nabla F(y_k) \quad \text{and} \quad y_{k+1} = x_{k+1} + \beta_{k+1}(x_{k+1} - x_k),$$

and, given (6), we obtain that

$$x_{k+1} = p(y_k) \quad \text{and} \quad y_{k+1} = x_{k+1} + \beta_{k+1}(x_{k+1} - x_k).$$

This is known as the accelerated proximal point algorithm introduced by Güler [27], which was recently extended in [33, 34].

Variable metric proximal point algorithm. Quasi-Newton methods can also be applied on F , which yields

$$(10) \quad x_{k+1} = x_k - \alpha_k B_k^{-1} \nabla F(x_k),$$

where B_k is the Hessian approximation of F based on quasi-Newton methods. This is known as the variable metric proximal point algorithm [7, 14, 22, 23].

Toward an inexact variable metric proximal point algorithm. Quasi-Newton approaches have been applied after the inexact Moreau envelope in various ways [7, 14, 22, 23]. In particular, it is shown in [14] that if the subproblems (5) are solved up to high enough accuracy, then the inexact variable metric proximal point algorithm preserves the superlinear convergence rate. However, the complexity of solving the subproblems with high accuracy is typically not taken into account in the above-mentioned works.

In the unrealistic case where $p(x_k)$ can be obtained at no cost, the proximal point algorithm can afford much larger step sizes than classical gradient methods, and thus

is more effective. For instance, when f is strongly convex, the Moreau envelope F can be made arbitrarily well conditioned by making κ arbitrarily small, according to (8). Then, a single gradient step on F is enough to be arbitrarily close to the optimum. In practice, however, subproblems are solved only approximately, and the complexity of solving the subproblems is directly related to the smoothing parameter κ . This leaves an important question: how large to choose the smoothing parameter κ . A small κ makes the smoothed function F better conditioned, while a large κ is needed to improve the conditioning of the subproblem (5).

The main contribution of our paper is to close this gap by providing a global complexity analysis that takes into account the complexity of solving the subproblems. More concretely, in the proposed QNing algorithm, we provide (i) a practical stopping criterion for the subproblems, (ii) several warm-start strategies, (iii) a simple line-search strategy that guarantees a sufficient descent in terms of function value. These three components together yield the global convergence analysis, which allows us to use the first-order method as a subproblem solver. Moreover, the global complexity we develop depends on the smoothing parameter κ , which provides some insight into how large to choose this parameter practically.

Solving the subproblems with first-order algorithms. In the composite setting, both proximal quasi-Newton methods and the variable metric proximal point algorithm require us to solve subproblems (4) and (5), respectively. In the general case, when a generic first-order method, e.g., proximal gradient descent, is used, our worst-case complexity analysis does not provide a clear winner, and our experiments in section 5.4 confirm that both approaches perform similarly. However, when it is possible to exploit the specific structure of the subproblems in one case but not in the other, the conclusion may differ.

For instance, when the problem has a finite sum (2) structure, the proximal point algorithm approach leads to subproblems that can be solved in $O(n \log(1/\varepsilon))$ iterations with first-order incremental methods such as SVRG [60], SAGA [15], or MISO [38], by using the same choice of smoothing parameter $\kappa = 2L/n$ as Catalyst [34]. Assuming that computing a gradient of a function f_i and computing the proximal operator of ψ are both feasible in $O(d)$ floating-point operations, our approach solves each subproblem with enough accuracy in $\tilde{O}(nd)$ operations.¹ On the other hand, we cannot naively apply SVRG to solve the proximal quasi-Newton update (4) at the same cost for the following reasons. First, the variable metric matrix B_k does not admit a natural finite sum decomposition. The naive way of writing it into n copies results in an increase in computational complexity for evaluating the incremental gradients. More precisely, when B_k has rank l , computing a single gradient now requires us to compute a matrix-vector product with cost at least $O(dl)$, resulting in an l -fold increase per iteration. Second, the previous iteration complexity $O(n \log(1/\varepsilon))$ for solving the subproblems would require the subproblems to be well conditioned, i.e., $B_k \succeq (L/n)I$, forcing the quasi-Newton metric to be potentially more isotropic. For these reasons, existing attempts to combine SVRG with quasi-Newton principles have taken other directions [26, 44].

3. QNing: A quasi-Newton meta-algorithm. We now present the QNing method in Algorithm 1, which consists of applying variable metric algorithms on the smoothed objective F with inexact gradients. Each gradient approximation is the result of a minimization problem tackled with the algorithm \mathcal{M} , used as a

¹The notation \tilde{O} hides logarithmic quantities.

subroutine. The outer loop of the algorithm performs quasi-Newton updates. The method \mathcal{M} can be any algorithm of the user's choice, as long as it enjoys a linear convergence rate for strongly convex problems. More technical details are given in section 3.1.

Algorithm 1 QNing: A quasi-Newton meta-algorithm.

input Initial point x_0 in \mathbb{R}^d ; number of iterations K ; smoothing parameter $\kappa > 0$; optimization algorithm \mathcal{M} ; optionally, budget $T_{\mathcal{M}}$ for solving the subproblems.

- 1: Initialization: $(g_0, F_0, z_0) = \text{ApproxGradient}(x_0, \mathcal{M})$; $H_0 = \frac{1}{\kappa}I$.
 - 2: **for** $k = 0, \dots, K - 1$ **do**
 - 3: Initialize $\eta_k = 1$.
 - 4: Perform the quasi-Newton step

$$x_{\text{test}} = x_k - (\eta_k H_k + (1 - \eta_k) H_0) g_k.$$
 - 5: Estimate the gradient and function value of the Moreau envelope at x_{test} :

$$(g_{\text{test}}, F_{\text{test}}, z_{\text{test}}) = \text{ApproxGradient}(x_{\text{test}}, \mathcal{M}).$$
 - 6: **while** $F_{\text{test}} > F_k - \frac{1}{4\kappa} \|g_k\|^2$ **do**
 - 7: Decrease the line-search parameter η_k in $[0, 1]$ and re-evaluate x_{test} .
 - 8: Re-estimate $(g_{\text{test}}, F_{\text{test}}, z_{\text{test}}) = \text{ApproxGradient}(x_{\text{test}}, \mathcal{M})$.
 - 9: **end while**
 - 10: Accept the new iterate: $(x_{k+1}, g_{k+1}, F_{k+1}, z_{k+1}) = (x_{\text{test}}, g_{\text{test}}, F_{\text{test}}, z_{\text{test}})$.
 - 11: Update H_{k+1} (for example, use L-BFGS update with $s_k = x_{k+1} - x_k$ and $y_k = g_{k+1} - g_k$).
 - 12: **end for**
- output** Inexact proximal point z_K (solution).
-

3.1. The main algorithm. We now discuss the main algorithm components and its main features.

Outer loop: Inexact variable metric proximal point algorithm. We apply variable metric algorithms with a simple line-search strategy similar to that of [54] on the Moreau envelope F . Given a positive definite matrix H_k and a step size η_k in $[0, 1]$, the algorithm computes the update

$$(LS) \quad x_{k+1} = x_k - (\eta_k H_k + (1 - \eta_k) H_0) g_k,$$

where $H_0 = (1/\kappa)I$. When $\eta_k = 1$, the update uses the metric H_k , and when $\eta_k = 0$, it uses an inexact proximal point update $x_{k+1} = x_k - (1/\kappa)g_k$. In other words, when the quality of the metric H_k is not good enough, due to the inexactness of the gradients used in its construction, the update is corrected toward a simple proximal point update whose convergence is well understood when the gradients are inexact.

In order to determine the step size η_k , we introduce the following descent condition:

$$(11) \quad F_{k+1} \leq F_k - \frac{1}{4\kappa} \|g_k\|^2.$$

We show that the descent condition (11) is always satisfied when $\eta_k = 0$; thus, the finite termination of the line search follows (see section 4.3 for more details). In our

experiments, we observed that the step size $\eta_k = 1$ was almost always selected. In practice, we try the values η_k in $\{1, 1/2, 1/4, 1/8, 0\}$ starting from the largest one and stopping when condition (11) is satisfied.

Example of variable metric algorithm: Inexact L-BFGS method. The L-BFGS rule we consider is the standard one and consists in updating incrementally a generating list of vectors $\{(s_i, y_i)\}_{i=1, \dots, j}$, which implicitly defines the L-BFGS matrix. We use here the two-loop recursion detailed in [50, Algorithm 7.4] and use skipping steps when the condition $s_i^\top y_i > 0$ is not satisfied, in order to ensure the positive-definiteness of the L-BFGS matrix H_k (see [20]).

Inner loop: Approximate the Moreau envelope. The inexactness of our scheme comes from the approximation of the Moreau envelope F and its gradient. The procedure `ApproxGradient`(\cdot) calls an minimization algorithm \mathcal{M} and applies \mathcal{M} to minimize the subproblem (14). When the problem is solved exactly, the function returns the exact values $g = \nabla F(x)$, $F_a = F(x)$, and $z = p(x)$. However, this is infeasible in practice and we can only expect approximate solutions. In particular, a stopping criterion should be specified. We consider the following variants.

- (a) We define an adaptive stopping criterion based on function values, and stop \mathcal{M} when the approximate solution satisfies the inequality (15). In contrast to the standard stopping criterion where the accuracy is an absolute constant, our stopping criterion is adaptive, since the right-hand side of (15) also depends on the current iterate z . More detailed theoretical insights will be given in section 4. Typically, checking whether or not the criterion is satisfied requires computing a duality gap, as in Catalyst [34].
- (b) We use a predefined budget $T_{\mathcal{M}}$ in terms of number of iterations of method \mathcal{M} , where $T_{\mathcal{M}}$ is a constant independent of k .

Note that such an adaptive stopping criterion is relatively classical in the literature of inexact gradient-based methods [8]. As we will see later, in section 4, when $T_{\mathcal{M}}$ is large enough, criterion (15) is guaranteed.

Requirements on \mathcal{M} . To apply QNing, the optimization method \mathcal{M} needs to have linear convergence rates for strongly convex problems. More precisely, for any strongly convex objective h , method \mathcal{M} should be able to generate a sequence of iterates $(w_t)_{t \geq 0}$ such that

$$(12) \quad h(w_t) - h^* \leq C_{\mathcal{M}}(1 - \tau_{\mathcal{M}})^t(h(w_0) - h^*)$$

for some constants $C_{\mathcal{M}} > 0$ and $1 > \tau_{\mathcal{M}} > 0$,

where w_0 is the initial point given to \mathcal{M} . The notion of linearly convergent methods extends naturally to nondeterministic methods where (12) is satisfied in expectation:

$$(13) \quad \mathbb{E}[h(w_t) - h^*] \leq C_{\mathcal{M}}(1 - \tau_{\mathcal{M}})^t(h(w_0) - h^*).$$

The linear convergence condition typically holds for many primal gradient-based optimization techniques, including classical full gradient descent methods, block-coordinate descent algorithms [47, 52], and variance-reduced incremental algorithms [15, 56, 60]. In particular, our method provides a generic way to combine incremental algorithms with quasi-Newton methods that are suitable for large-scale optimization problems. For simplicity, we only consider the deterministic variant (12) in the analysis. However, it is possible to show that the same complexity results still hold for nondeterministic methods in expectation, as discussed in section 4.5. We emphasize

that we do not assume any convergence guarantee of \mathcal{M} on nonstrongly convex problems, since our subproblems are always strongly convex.

Warm starts for the subproblems. The employment of an adequate initialization for solving each subproblem plays an important role in our analysis. The warm-start strategy we propose here ensures that the stopping criterion in each subproblem can be achieved in a constant number of iterations.

Consider the minimization of a subproblem

$$\min_{w \in \mathbb{R}^d} \left\{ h(w) \triangleq f(w) + \frac{\kappa}{2} \|w - x\|^2 \right\}.$$

Then, our warm-start strategy depends on the nature of f :

- when f is smooth, we initialize with $w_0 = x$;
- when $f = f_0 + \psi$ is composite, we initialize with

$$w_0 = \arg \min_{w \in \mathbb{R}^d} \left\{ f_0(x) + \langle \nabla f_0(x), w - x \rangle + \frac{L + \kappa}{2} \|w - x\|^2 + \psi(w) \right\},$$

which performs an additional proximal step compared to the smooth case.

Handling composite objective functions. In machine learning or signal processing, convex composite objectives (1) with a nonsmooth penalty ψ are typically formulated to encourage solutions with specific characteristics; in particular, the ℓ_1 -norm is known to provide sparsity. Smoothing techniques [46] may allow us to solve the optimization problem up to some chosen accuracy, but they provide solutions that do not inherit the properties induced by the nonsmoothness of the objective. To illustrate what we mean by this statement, we may consider smoothing the ℓ_1 -norm, leading to a solution vector with small coefficients but not exact zeros. When the goal is to perform model selection, that is, understanding which variables are important to explain a phenomenon, exact sparsity is seen as an asset, and optimization techniques dedicated to composite problems such as the fast iterative shrinkage-thresholding algorithm (FISTA) [2] are often preferred (see [40]).

One might then be concerned that our scheme operates on the smoothed objective F , leading to iterates $(x_k)_{k \geq 0}$ that may suffer from the above “nonsparse” issue, assuming that ψ is the ℓ_1 -norm. Yet, our approach does not directly output the iterates $(x_k)_{k \geq 0}$ but rather their proximal mappings $(z_k)_{k \geq 0}$. In particular, the ℓ_1 -regularization is encoded in the proximal mapping (14). Thus, the approximate proximal point z_k may be sparse. For this reason, our theoretical analysis presented in section 4 studies the convergence of the sequence $(f(z_k))_{k \geq 0}$ to the solution f^* .

4. Convergence and complexity analysis. In this section, we study the convergence of the QNing algorithm, that is, the rate of convergence of the quantities $(F(x_k) - F^*)_{k \geq 0}$ and $(f(z_k) - f^*)_{k \geq 0}$, and also the computational complexity due to solving the subproblems (14). We start by stating the main properties of the gradient approximation in section 4.1. Then, we analyze the convergence of the outer loop algorithm in section 4.2, and section 4.3 is devoted to the properties of the line-search strategy. After that, we provide the cost of solving the subproblems in section 4.4 and derive the global complexity analysis in section 4.5.

4.1. Properties of the gradient approximation. The next lemma is classical and provides approximation guarantees about the quantities returned by the ApproxGradient procedure (Algorithm 2); see [5, 23]. We recall here the proof for completeness.

Algorithm 2 Generic procedure ApproxGradient.

input Current point x in \mathbb{R}^d ; smoothing parameter $\kappa > 0$; optionally, budget $T_{\mathcal{M}}$.

1: Compute the approximate proximal mapping using an optimization method \mathcal{M} :

$$(14) \quad z \approx \arg \min_{w \in \mathbb{R}^d} \left\{ h(w) \triangleq f(w) + \frac{\kappa}{2} \|w - x\|^2 \right\},$$

using one of the following stopping criteria:

- stop when the approximate solution z satisfies

$$(15) \quad h(z) - h^* \leq \frac{\kappa}{36} \|z - x\|^2;$$

- stop when we reach the predefined constant budget $T_{\mathcal{M}}$ (for instance, one pass over the data).

2: Estimate the gradient $\nabla F(x)$ of the Moreau envelope using

$$g = \kappa(x - z).$$

output Gradient estimate g , objective value estimate $F_a \triangleq h(z)$, proximal point estimate z .

LEMMA 1 (approximation quality of the gradient approximation). *Consider a point x in \mathbb{R}^d , a positive scalar ε , and an approximate proximal point*

$$z \approx \arg \min_{w \in \mathbb{R}^d} \left\{ h(w) \triangleq f(w) + \frac{\kappa}{2} \|w - x\|^2 \right\},$$

such that

$$h(z) - h^* \leq \varepsilon,$$

where $h^* = \min_{w \in \mathbb{R}^d} h(w)$. As in Algorithm 2, we define the gradient estimate $g = \kappa(x - z)$ and the function value estimate $F_a = h(z)$. Then, the following inequalities hold:

$$(16) \quad F(x) \leq F_a \leq F(x) + \varepsilon,$$

$$(17) \quad \|z - p(x)\| \leq \sqrt{\frac{2\varepsilon}{\kappa}},$$

$$(18) \quad \|g - \nabla F(x)\| \leq \sqrt{2\kappa\varepsilon}.$$

Moreover, F_a is related to f by the following relationship:

$$(19) \quad f(z) = F_a - \frac{1}{2\kappa} \|g\|^2.$$

Proof. Relations (16) and (19) are straightforward by the definition of $h(z)$. Since f is convex, the function h is κ -strongly convex, and (17) follows from

$$\frac{\kappa}{2} \|z - p(x)\|^2 \leq h(z) - h(p(x)) = h(z) - h^* \leq \varepsilon,$$

where we recall that $p(x)$ minimizes h . Finally, we obtain (18) from

$$g - \nabla F(x) = \kappa(x - z) - \kappa(x - p(x)) = \kappa(p(x) - z),$$

by using the definitions of g and the property (6). \square

This lemma allows us to quantify the quality of the gradient and function value approximations, which is crucial to control the error accumulation of inexact proximal point methods. Moreover, the relation (19) establishes a link between the approximate function value of F and the function value of the original objective f ; as a consequence, it is possible to relate the convergence rate of f to the convergence rate of F . Finally, the following result is a direct consequence of Lemma 1.

LEMMA 2 (bounding the exact gradient by its approximation). *Consider the quantities introduced in Lemma 1. Then,*

$$(20) \quad \frac{1}{2}\|g\|^2 - 2\kappa\varepsilon \leq \|\nabla F(x)\|^2 \leq 2(\|g\|^2 + 2\kappa\varepsilon).$$

Proof. The right-hand side of (20) follows from

$$\begin{aligned} \|\nabla F(x)\|^2 &\leq 2(\|\nabla F(x) - g\|^2 + \|g\|^2) \\ &\leq 2(2\kappa\varepsilon + \|g\|^2) \quad (\text{from (18)}). \end{aligned}$$

Interchanging $\nabla F(x)$ and g gives the left-hand side of the inequality. \square

COROLLARY 1. *If $\varepsilon \leq \frac{c}{\kappa}\|g\|^2$ with $c < \frac{1}{4}$, then*

$$(21) \quad \frac{1-4c}{2} \leq \frac{\|\nabla F(x)\|^2}{\|g\|^2} \leq 2(1+2c).$$

This corollary is important since it allows us to replace the unknown exact gradient $\|\nabla F(x)\|$ by its approximation $\|g\|$, at the cost of a constant factor, as long as the condition $\varepsilon \leq \frac{c}{\kappa}\|g\|^2$ is satisfied.

4.2. Convergence analysis of the outer loop. We are now in a position to establish the convergence of the QNing meta-algorithm, without yet considering the cost of solving the subproblems (14). At iteration k , an approximate proximal point is evaluated:

$$(22) \quad (g_k, F_k, z_k) = \text{ApproxGradient}(x_k, \mathcal{M}).$$

The following lemma characterizes the expected descent in terms of objective function value.

LEMMA 3 (approximate descent property). *At iteration k , if the subproblem (22) is solved up to accuracy ε_k in the sense of Lemma 1 and the next iterate x_{k+1} satisfies the descent condition (11), then*

$$(23) \quad F(x_{k+1}) \leq F(x_k) - \frac{1}{8\kappa}\|\nabla F(x_k)\|^2 + \frac{3}{2}\varepsilon_k.$$

Proof. From (16) and (11),

$$\begin{aligned} F(x_{k+1}) &\leq F_{k+1} \leq F_k - \frac{1}{4\kappa}\|g_k\|^2 \\ &\leq F(x_k) + \varepsilon_k - \left(\frac{1}{8\kappa}\|\nabla F(x_k)\|^2 - \frac{\varepsilon_k}{2} \right) \quad (\text{from (16) and (20)}) \\ &= F(x_k) - \frac{1}{8\kappa}\|\nabla F(x_k)\|^2 + \frac{3}{2}\varepsilon_k. \end{aligned} \quad \square$$

This lemma gives us an initial clue about the natural choice of the accuracy ε_k , which should be of the same order as $\|\nabla F(x_k)\|^2$. In particular, if

$$(24) \quad \varepsilon_k \leq \frac{1}{16\kappa} \|\nabla F(x_k)\|^2,$$

then we have

$$(25) \quad F(x_{k+1}) \leq F(x_k) - \frac{1}{32\kappa} \|\nabla F(x_k)\|^2,$$

which is a typical inequality used for analyzing gradient descent methods. Before presenting the convergence result, we remark that condition (24) cannot be used directly since it requires the exact gradient $\|\nabla F(x_k)\|$ to be known. A more practical choice consists of replacing it by the approximate gradient.

LEMMA 4 (a practical choice of ε_k). *The following condition implies inequality (24):*

$$(26) \quad \varepsilon_k \leq \frac{1}{36\kappa} \|g_k\|^2.$$

Proof. From Corollary 1, (26) implies

$$\|g_k\|^2 \leq \frac{2}{1 - \frac{4}{36}} \|\nabla F(x_k)\|^2 = \frac{9}{4} \|\nabla F(x_k)\|^2,$$

and thus

$$\varepsilon_k \leq \frac{1}{36\kappa} \|g_k\|^2 \leq \frac{1}{16\kappa} \|\nabla F(x_k)\|^2. \quad \square$$

This is the first stopping criterion (15) in Algorithm 2. Finally, we obtain the following convergence result for strongly convex problems, which is classical in the literature of inexact gradient methods (see [8, section 4.1] for a similar result).

PROPOSITION 2 (convergence of Algorithm 1, strongly convex objectives). *Assume that f is μ -strongly convex. Let $(x_k)_{k \geq 0}$ be the sequences generated by Algorithm 1 where the stopping criterion (15) is used. Then,*

$$F(x_k) - F^* \leq \left(1 - \frac{1}{16q}\right)^k (F(x_0) - F^*), \quad \text{with } q = \frac{\mu + \kappa}{\mu}.$$

Proof. The proof follows directly from (25) and the standard analysis of the GD algorithm for the μ_F -strongly convex and L_F -smooth function F by remarking that $L_F = \kappa$ and $\mu_F = \frac{\mu\kappa}{\mu + \kappa}$. \square

COROLLARY 2. *Under the conditions of Proposition 2, we have*

$$(27) \quad f(z_k) - f^* \leq \left(1 - \frac{1}{16q}\right)^k (f(x_0) - f^*).$$

Proof. From (19) and (26), we have

$$f(z_k) = F_k - \frac{1}{2\kappa} \|g_k\|^2 \leq F(x_k) + \varepsilon_k - \frac{1}{2\kappa} \|g_k\|^2 \leq F(x_k).$$

Moreover, $F(x_0)$ is upper-bounded by $f(x_0)$, following (7). \square

It is worth pointing out that our analysis establishes a linear convergence rate, whereas one would expect a superlinear convergence rate as for classical variable metric methods. The trade-off lies in the choice of the approximation error ε_k . In order to achieve superlinear convergence, the approximation error ε_k needs to decrease superlinearly, as shown in [14]. However, a quickly decreasing sequence ε_k requires increasing effort in solving the subproblems, which will dominate the global complexity. In other words, the global complexity may increase even though we achieve faster convergence in the outer loop. This will become clearer when we discuss the inner loop complexity in section 4.4.

Next, we show that, under a bounded level set condition, QNing enjoys the classical sublinear $O(1/k)$ convergence rate when the objective is convex but not strongly convex.

PROPOSITION 3 (convergence of Algorithm 1 for convex but not strongly convex objectives). *Let f be a convex function with bounded level sets. Then, there exists a constant $R > 0$ that depends on the initialization point x_0 , such that the sequences $(x_k)_{k \geq 0}$ and $(z_k)_{k \geq 0}$ generated by Algorithm 1 with stopping criterion (15) satisfy*

$$F(x_k) - F^* \leq \frac{32\kappa R^2}{k} \quad \text{and} \quad f(z_k) - f^* \leq \frac{32\kappa R^2}{k}.$$

Proof. We defer the proof and the proper definition of the bounded level set assumption to Appendix A. \square

So far, we have assumed in our analysis that the iterates satisfy the descent condition (11), which means the line-search strategy will always terminate. We prove in the next section that this is indeed the case, and provide some additional conditions under which a nonzero step size will be selected.

4.3. Conditions for nonzero step sizes η_k and termination of the line search. At iteration k , a line search is performed on the step size η_k to find the next iterate

$$x_{k+1} = x_k - (\eta_k H_k + (1 - \eta_k) H_0) g_k,$$

such that x_{k+1} satisfies the descent condition (11). We first show that the descent condition holds when $\eta_k = 0$, before giving a more general result.

LEMMA 5. *If the subproblems are solved up to accuracy $\varepsilon_k \leq \frac{1}{36\kappa} \|g_k\|^2$, then the descent condition (11) holds when $\eta_k = 0$.*

Proof. When $\eta_k = 0$, $x_{k+1} = x_k - \frac{1}{\kappa} g_k = z_k$. Then,

$$\begin{aligned} F_{k+1} &\leq F(x_{k+1}) + \frac{1}{36\kappa} \|g_{k+1}\|^2 \quad (\text{from (16)}) \\ &\leq F(x_{k+1}) + \frac{1}{36\kappa} \frac{2}{1 - \frac{4}{36}} \|\nabla F(x_{k+1})\|^2 \quad (\text{from (21)}) \\ &< F(x_{k+1}) + \frac{1}{2\kappa} \|\nabla F(z_{k+1})\|^2 \\ &\leq f(x_{k+1}) = f(z_k) \quad (\text{from (7)}) \\ &= F_k - \frac{1}{2\kappa} \|g_k\|^2 \quad (\text{from (19)}). \end{aligned}$$

\square

Therefore, it is theoretically sound to take the trivial step size $\eta_k = 0$, which implies the termination of our line-search strategy. In other words, the descent condition always holds by taking an inexact gradient step on the Moreau envelope F , which corresponds to the update of the proximal point algorithm. However, the purpose of using the variable metric method is to exploit the curvature of the function, which is not the case when $\eta_k = 0$. Thus, the trivial step size should only be considered as a backup plan, and we show in the following some sufficient conditions for taking nonzero step sizes and even stronger conditions for unit step sizes.

LEMMA 6 (a sufficient condition for satisfying the descent condition (11)). *If the subproblems are solved up to accuracy $\varepsilon_k \leq \frac{1}{36\kappa}\|g_k\|^2$, then the sufficient condition (11) holds for any $x_{k+1} = x_k - A_k g_k$, where A_k is a positive definite matrix satisfying $\frac{1-\alpha}{\kappa}I \preceq A_k \preceq \frac{1+\alpha}{\kappa}I$ with $\alpha \leq \frac{1}{3}$.*

As a consequence, a line-search strategy consisting of finding the largest η_k of the form γ^i , with $i = 1, \dots, +\infty$ and γ in $(0, 1)$, always terminates in a bounded number of iterations if the sequence of variable metrics $(H_k)_{k \geq 0}$ is bounded, i.e., there exists $0 < m < M$ such that, for any k , $mI \preceq H_k \preceq MI$. This is the case for the L-BFGS update.

LEMMA 7 (boundedness of L-BFGS metric matrix [50, Chapters 8 and 9]). *The variable metric matrices $(B_k)_k$ constructed by the L-BFGS rule are positive definite and bounded.*

Proof of Lemma 6. First, we recall that $z_k = x_k - \frac{1}{\kappa}g_k$ and we rewrite

$$F_{k+1} = \underbrace{F_{k+1} - F(x_{k+1})}_{\triangleq E_1} + \underbrace{F(x_{k+1}) - F(z_k)}_{\triangleq E_2} + F(z_k)$$

as follows. We shall bound the two error terms E_1 and E_2 by some factors of $\|g_k\|^2$. Noting that the subproblems are solved up to $\varepsilon_k \leq \frac{c}{\kappa}\|g_k\|^2$ with $c = \frac{1}{36}$, we obtain by construction that

$$(28) \quad E_1 = F_{k+1} - F(x_{k+1}) \leq \varepsilon_{k+1} \leq \frac{c}{\kappa}\|g_{k+1}\|^2 \leq \frac{2c}{(1-4c)\kappa}\|\nabla F(x_{k+1})\|^2,$$

where the last inequality comes from Corollary 1. Moreover,

$$\begin{aligned} \|\nabla F(x_{k+1})\| &\leq \|\nabla F(z_k)\| + \|\nabla F(x_{k+1}) - \nabla F(z_k)\| \\ &\leq \|\nabla F(z_k)\| + \kappa\|x_{k+1} - z_k\|. \end{aligned}$$

Since $x_{k+1} - z_k = (\frac{1}{\kappa} - A_k)g_k$, we have $\|x_{k+1} - z_k\| \leq \frac{\alpha}{\kappa}\|g_k\|$. This implies that

$$(29) \quad \|\nabla F(x_{k+1})\| \leq \|\nabla F(z_k)\| + \alpha\|g_k\|,$$

and thus

$$(30) \quad E_1 \leq \frac{4c}{(1-4c)\kappa} (\|\nabla F(z_k)\|^2 + \alpha^2\|g_k\|^2).$$

Second, by the κ -smoothness of F , we have

$$\begin{aligned}
 E_2 &= F(x_{k+1}) - F(z_k) \\
 &\leq \langle \nabla F(z_k), x_{k+1} - z_k \rangle + \frac{\kappa}{2} \|x_{k+1} - z_k\|^2 \\
 &\leq \frac{1}{4\kappa} \|\nabla F(z_k)\|^2 + \kappa \|x_{k+1} - z_k\|^2 + \frac{\kappa}{2} \|x_{k+1} - z_k\|^2 \\
 (31) \quad &\leq \frac{1}{4\kappa} \|\nabla F(z_k)\|^2 + \frac{3\alpha^2}{2\kappa} \|g_k\|^2,
 \end{aligned}$$

where the last inequality follows from $\|x_{k+1} - z_k\| \leq \frac{\alpha}{\kappa} \|g_k\|$. Combining (30) and (31) yields

$$(32) \quad E_1 + E_2 \leq \left[\frac{4c}{1-4c} + \frac{1}{4} \right] \frac{1}{\kappa} \|\nabla F(z_k)\|^2 + \left[\frac{4c}{1-4c} + \frac{3}{2} \right] \frac{\alpha^2}{\kappa} \|g_k\|^2.$$

When $c \leq \frac{1}{36}$ and $\alpha \leq \frac{1}{3}$, we have

$$E_1 + E_2 \leq \frac{1}{2\kappa} \|\nabla F(z_k)\|^2 + \frac{1}{4\kappa} \|g_k\|^2.$$

Therefore,

$$\begin{aligned}
 F_{k+1} &\leq F(z_k) + E_1 + E_2 \\
 &\leq F(z_k) + \frac{1}{2\kappa} \|\nabla F(z_k)\|^2 + \frac{1}{4\kappa} \|g_k\|^2 \\
 &\leq f(z_k) + \frac{1}{4\kappa} \|g_k\|^2 \\
 (33) \quad &= F_k - \frac{1}{4\kappa} \|g_k\|^2,
 \end{aligned}$$

where the last equality follows from (19). This completes the proof. \square

Note that, in practice, we consider a set of step sizes $\eta_k = \gamma^i$ for $i \leq i_{\max}$ or $\eta_k = 0$, which naturally upper-bounds the number of line-search iterations to i_{\max} . More precisely, all experiments performed in this paper use $\gamma = 1/2$ and $i_{\max} = 3$. Moreover, we observe that the unit step size is very often sufficient for the descent condition to hold, as studied empirically in Appendix C.2.

The following result shows that, under a specific assumption on the Moreau envelope F , the unit step size is indeed selected when the iterate are close to the optimum. The condition, called the Dennis–Moré criterion [17], is classical in the literature of quasi-Newton methods. Even though we cannot formally show that the criterion holds for the Moreau envelope F , since it requires F to be twice continuously differentiable, which is not true in general (see [32]), it provides a sufficient condition for the unit step size. Therefore, the lemma below should be seen not as a formal explanation for the choice of step size $\eta_k = 1$, but simply as a reasonable condition that leads to this choice.

LEMMA 8 (a sufficient condition for unit step size). *Assume that f is strongly convex and F is twice continuously differentiable with Lipschitz continuous Hessian $\nabla^2 F$. If the subproblems are solved up to accuracy $\varepsilon_k \leq \frac{\mu^2}{128\kappa(\mu+\kappa)^2} \|g_k\|^2$ and the Dennis–Moré criterion [17] is satisfied, i.e.,*

$$(DM) \quad \lim_{k \rightarrow \infty} \frac{\|(B_k^{-1} - \nabla^2 F(x^*)^{-1})g_k\|}{\|g_k\|} = 0,$$

where x^* is the minimizer of the problem and $B_k = H_k^{-1}$ is the variable metric matrix, then the descent condition (11) is satisfied with $\eta_k = 1$ when k is large enough.

We remark that the Dennis–Moré criterion we use here is slightly different from the standard one since our criterion is based on approximate gradients g_k . If the g_k 's are indeed the exact gradients and the variable metrics B_k are bounded, then our criterion is equivalent to the standard Dennis–Moré criterion. The proof of the lemma is close to that of similar lemmas appearing in the proximal quasi-Newton literature [31], and is relegated to the appendix. Interestingly, this proof also suggests that a stronger stopping criterion ε_k such that $\varepsilon_k = o(\|g_k\|^2)$ could lead to superlinear convergence. However, such a choice of ε_k would significantly increase the complexity for solving the subproblems, and degrade the overall complexity.

4.4. Complexity analysis of the inner loop. In this section, we evaluate the complexity of solving the subproblems (14) up to the desired accuracy using a linearly convergent method \mathcal{M} . Our main result is that all subproblems can be solved in a constant number $T_{\mathcal{M}}$ of iterations (in expectation if the method is nondeterministic) using the proposed warm-start strategy. Let us consider the subproblem with an arbitrary prox center x ,

$$(34) \quad \min_{w \in \mathbb{R}^d} \left\{ h(w) = f(w) + \frac{\kappa}{2} \|w - x\|^2 \right\}.$$

The number of iterations needed is determined by the ratio between the initial gap $h(w_0) - h^*$ and the desired accuracy. We shall bound this ratio by a constant factor.

LEMMA 9 (warm start for primal methods: smooth case). *If f is differentiable with L -Lipschitz continuous gradients, we initialize the method \mathcal{M} with $w_0 = x$. Then, we have the guarantee that*

$$(35) \quad h(w_0) - h^* \leq \frac{L + \kappa}{2\kappa^2} \|\nabla F(x)\|^2.$$

Proof. Denote by w^* the minimizer of h . Then, we have the optimality condition $\nabla f(w^*) + \kappa(w^* - x) = 0$. As a result,

$$\begin{aligned} h(w_0) - h^* &= f(x) - \left(f(w^*) + \frac{\kappa}{2} \|w^* - x\|^2 \right) \\ &\leq f(w^*) + \langle \nabla f(w^*), x - w^* \rangle + \frac{L}{2} \|x - w^*\|^2 - \left(f(w^*) + \frac{\kappa}{2} \|w^* - x\|^2 \right) \\ &= \frac{L + \kappa}{2} \|w^* - x\|^2 \\ &= \frac{L + \kappa}{2\kappa^2} \|\nabla F(x)\|^2. \quad \square \end{aligned}$$

The inequality in the above proof relies on the smoothness of f , which does not hold for composite problems. The next lemma addresses this issue.

LEMMA 10 (warm start for primal methods: composite case). *Consider the composite optimization problem $f = f_0 + \psi$, where f_0 is L -smooth. By initializing with*

$$(36) \quad w_0 = \arg \min_{w \in \mathbb{R}^d} \left\{ f_0(x) + \langle \nabla f_0(x), w - x \rangle + \frac{L + \kappa}{2} \|w - x\|^2 + \psi(w) \right\},$$

we have

$$h(w_0) - h^* \leq \frac{L + \kappa}{2\kappa^2} \|\nabla F(x)\|^2.$$

Proof. We use the inequality corresponding to [2, Lemma 2.3]: for any w ,

$$(37) \quad h(w) - h(w_0) \geq \frac{L'}{2} \|w_0 - x\|^2 + L' \langle w_0 - x, x - w \rangle,$$

with $L' = L + \kappa$. Then, we apply this inequality to $w = w^*$ to obtain

$$\begin{aligned} h(w_0) - h^* &\leq -\frac{L'}{2} \|w_0 - x\|^2 - L' \langle w_0 - x, x - w^* \rangle \\ &\leq \frac{L'}{2} \|x - w^*\|^2 = \frac{L + \kappa}{2\kappa^2} \|\nabla F(x)\|^2. \end{aligned} \quad \square$$

We get an initialization in the composite case of the same quality as that in the smooth case by performing an additional proximal step. It is important to remark that the above analysis does not require the strong convexity of f , which allows us to derive the desired inner-loop complexity.

PROPOSITION 4 (inner-loop complexity for Algorithm 1). *Consider Algorithm 1 with the warm-start strategy described in either Lemma 9 or Lemma 10. Assume that the optimization method \mathcal{M} applied in the inner loop produces a sequence $(w_t)_{t \geq 0}$ for each subproblem (34) such that*

$$(38) \quad h(w_t) - h^* \leq C_{\mathcal{M}}(1 - \tau_{\mathcal{M}})^t (h(w_0) - h^*) \quad \text{for some constants } C_{\mathcal{M}}, \tau_{\mathcal{M}} > 0.$$

Then, the stopping criterion $\varepsilon \leq \frac{1}{72\kappa} \|g\|^2$ is achieved in at most $T_{\mathcal{M}}$ iterations with

$$T_{\mathcal{M}} = \frac{1}{\tau_{\mathcal{M}}} \log \left(38C_{\mathcal{M}} \frac{L + \kappa}{\kappa} \right).$$

Proof. Consider that at iteration k we apply \mathcal{M} to approximate the proximal mapping according to x . With the given $T_{\mathcal{M}}$ (which we abbreviate as T), we have

$$\begin{aligned} h(w_T) - h^* &\leq C_{\mathcal{M}}(1 - \tau_{\mathcal{M}})^T (h(w_0) - h^*) \\ &\leq C_{\mathcal{M}} e^{-\tau_{\mathcal{M}} T} (h(w_0) - h^*) \\ &\leq C_{\mathcal{M}} e^{-\tau_{\mathcal{M}} T} \frac{L + \kappa}{2\kappa^2} \|\nabla F(x)\|^2 \quad (\text{by Lemmas 9 and 10}) \\ &= \frac{1}{76\kappa} \|\nabla F(x)\|^2 \\ &\leq \frac{1}{36\kappa} \|g\|^2, \end{aligned}$$

where the last inequality follows from Lemma 2. \square

Next, we extend the previous result obtained with deterministic methods \mathcal{M} to randomized ones, where linear convergence is only achieved in the expectation. The proof is a simple application of [33, Lemma C.1] (see also [12] for related results on the expected complexity of randomized algorithms).

Remark 1 (\mathcal{M} is nondeterministic). Assume that the optimization method \mathcal{M} applied to each subproblem (34) produces a sequence $(w_t)_{t \geq 0}$ such that

$$\mathbb{E}[h(w_t) - h^*] \leq C_{\mathcal{M}}(1 - \tau_{\mathcal{M}})^t (h(w_0) - h^*) \quad \text{for some constants } C_{\mathcal{M}}, \tau_{\mathcal{M}} > 0.$$

We define the stopping time $T_{\mathcal{M}}$ by

$$(39) \quad T_{\mathcal{M}} = \inf \left\{ t \geq 1 \mid h(w_t) - h^* \leq \frac{1}{36\kappa} \|g_t\|^2 \right\}, \quad \text{where } g_t = \kappa(x - w_t),$$

which is the random variable corresponding to the minimum number of iterations to guarantee the stopping condition (15). Then, when the warm-start strategy described in either Lemma 9 or Lemma 10 is applied, the expected number of iterations satisfies

$$(40) \quad \mathbb{E}[T_{\mathcal{M}}] \leq \frac{1}{\tau_{\mathcal{M}}} \log \left(76C_{\mathcal{M}} \frac{L + \kappa}{\tau_{\mathcal{M}}\kappa} \right) + 1.$$

Remark 2 (checking the stopping criterion). Notice that the stopping criterion (15), i.e., $h(w) - h^* \leq \frac{\kappa}{36}\|w - x\|^2$, cannot be checked directly since h^* is unknown. Nevertheless, an upper bound on the optimality gap $h(w) - h^*$ is usually available. In particular,

- when f is smooth, which implies h is smooth, we have

$$(41) \quad h(w) - h^* \leq \frac{1}{2(\mu + \kappa)} \|\nabla h(w)\|^2,$$

- otherwise, we can evaluate the Fenchel conjugate function, which is a natural lower bound of h^* ; see [37, section D.2.3].

4.5. Global complexity of QNing. Finally, we can use the previous results to upper-bound the complexity of the QNing algorithm in terms of iterations of the method \mathcal{M} for minimizing f up to ε .

PROPOSITION 5 (worst-case global complexity for Algorithm 1). *Given a linearly convergent method \mathcal{M} satisfying (12), we apply \mathcal{M} to solve the subproblems of Algorithm 1 with the warm-start strategy given in either Lemma 9 or Lemma 10 up to accuracy $\varepsilon_k \leq \frac{1}{36\kappa}\|g_k\|^2$. Then, the number of iterations of method \mathcal{M} to guarantee the optimality condition $f(z_k) - f^* \leq \varepsilon$ is as follows:*

- for μ -strongly convex problems,

$$\begin{aligned} O \left(T_{\mathcal{M}} \times \frac{\mu + \kappa}{\mu} \log \left(\frac{f(x_0) - f^*}{\varepsilon} \right) \right) \\ = O \left(\frac{\mu + \kappa}{\tau_{\mathcal{M}}\mu} \log \left(\frac{f(x_0) - f^*}{\varepsilon} \right) \log \left(38C_{\mathcal{M}} \frac{L + \kappa}{\kappa} \right) \right); \end{aligned}$$

- for convex problems with bounded level sets,

$$O \left(T_{\mathcal{M}} \times \frac{2\kappa R^2}{\varepsilon} \right) = O \left(\frac{2\kappa R^2}{\tau_{\mathcal{M}}\varepsilon} \log \left(38C_{\mathcal{M}} \frac{L + \kappa}{\kappa} \right) \right).$$

Proof. The total number of calls of method \mathcal{M} is simply $T_{\mathcal{M}}$ times the number of outer-loop iterations times the potential number of line-search steps at each iteration (which is hidden in the $O(\cdot)$ notation since this number can be made arbitrarily small). \square

Remark 3. For nondeterministic methods, applying (40) yields a global complexity in expectation similar to the previous result with additional constant $2/\tau_{\mathcal{M}}$ in the last log factor.

As we shall see, the global complexity of our algorithm is mainly controlled by the smoothing parameter κ . Unfortunately, under the current analysis, our algorithm QNing does not lead to an improved convergence rate in terms of the worst-case complexity bounds. It is worthwhile underlining, however, that this result is not surprising since it is often the case for L-BFGS-type methods, for which there remains

an important gap between theory and practice. Indeed, L-BFGS often outperforms the vanilla GD method in many practical cases, but never in theory, which turns out to be the bottleneck in our analysis.

We give below the worst-case global complexity of QNing when applied to two optimization methods \mathcal{M} of interest. Proposition 5 and its application to the two examples show that, in terms of worst-case complexity, the QNing scheme leaves the convergence rate almost unchanged.

Example 1. Consider GD with fixed constant step size $1/L$ as the optimization method \mathcal{M} . Directly applying GD to minimize f requires

$$O(L/\mu \log(1/\varepsilon))$$

iterations to achieve ε accuracy. The complexity to achieve the same accuracy with QNing-GD is, in the worst case,

$$\tilde{O}((L + \kappa)/\mu \log(1/\varepsilon)).$$

Example 2. Consider the stochastic variance-reduced gradient (SVRG) as the optimization method \mathcal{M} . SVRG minimizes f to ε accuracy in

$$O\left(\max\left\{n, \frac{L}{\mu}\right\} \log\left(\frac{1}{\varepsilon}\right)\right)$$

iterations in expectation. QNing-SVRG achieves the same result with the worst-case expected complexity

$$\tilde{O}\left(\max\left\{\frac{\mu + \kappa}{\mu}n, \frac{L + \kappa}{\mu}\right\} \log\left(\frac{1}{\varepsilon}\right)\right).$$

Choice of κ . Minimizing the above worst-case complexity with respect to κ suggests that κ should be chosen as small as possible. However, such a statement is based on the pessimistic theoretical analysis of the L-BFGS-type method, which is not better than standard gradient descent methods. Noting that for smooth functions the L-BFGS method often outperforms Nesterov's accelerated gradient method, it is reasonable to expect they achieve a similar complexity bound. In other words, the choice of κ may be substantially different if one is able to show that the L-BFGS-type method enjoys an accelerated convergence rate.

In order to illustrate the difference, we heuristically assume that the L-BFGS method enjoys a similar convergence rate to Nesterov's accelerated gradient method. Then, the global complexity of our algorithm QNing matches the complexity of the related Catalyst acceleration scheme [33], which will be

$$\tilde{O}\left(\frac{1}{\tau_{\mathcal{M}}} \sqrt{\frac{\mu + \kappa}{\mu}} \log\left(\frac{1}{\varepsilon}\right)\right),$$

for μ -strongly convex problems. In such a case, the complexity of QNing-GD and QNing-SVRG will be

$$\tilde{O}\left(\frac{L + \kappa}{\sqrt{(\mu + \kappa)\mu}} \log\left(\frac{1}{\varepsilon}\right)\right) \quad \text{and} \quad \tilde{O}\left(\max\left\{\sqrt{\frac{\mu + \kappa}{\mu}}n, \frac{L + \kappa}{\sqrt{(\mu + \kappa)\mu}}\right\} \log\left(\frac{1}{\varepsilon}\right)\right),$$

which enjoy acceleration by taking $\kappa = O(L)$ and $\kappa = O(L/n)$, respectively. In the following section, we will experiment with this heuristic as if the L-BFGS method

enjoys an accelerated convergence rate. More precisely, we will choose the smoothing parameter κ as in the related Catalyst acceleration scheme [33], and present empirical evidence in support of this heuristic.

5. Experiments and practical details. In this section, we present the experimental results obtained by applying QNing to several first-order optimization algorithms. We start by presenting various benchmarks and practical parameter-tuning choices. Then, we study the performance of QNing applied to SVRG (section 5.3) and to the proximal gradient algorithm ISTA (iterative shrinkage-thresholding algorithm; see section 5.4), which reduces to GD in the smooth case. We demonstrate that QNing can be viewed as an acceleration scheme: by applying QNing to an optimization algorithm \mathcal{M} , we achieve better performance than when applying \mathcal{M} directly to the problem. In addition, we compare QNing to existing stochastic variants of the L-BFGS algorithm in section 5.3. Finally, we study the behavior of QNing under different choices of parameters in section 5.5. The code used for all the experiments is available at <https://github.com/hongzhoulin89/Catalyst-QNing/>.

5.1. Formulations and data sets. We consider three common optimization problems in machine learning and signal processing: logistic regression, least absolute shrinkage and selection operator (Lasso) regression, and linear regression with elastic-net regularization. These three formulations all admit the composite finite-sum structure but differ in terms of smoothness and strength of convexity. The three specific formulations are listed below.

- ℓ_2^2 -regularized logistic regression:

$$\min_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-b_i a_i^T x)) + \frac{\mu}{2} \|x\|^2,$$

which leads to a μ -strongly convex smooth optimization problem.

- ℓ_1 -regularized linear regression (Lasso):

$$\min_{x \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (b_i - a_i^T x)^2 + \lambda \|x\|_1,$$

which is convex and nonsmooth, but not strongly convex.

- $(\ell_1 - \ell_2^2)$ -regularized linear regression (elastic-net):

$$\min_{x \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (b_i - a_i^T x)^2 + \lambda \|x\|_1 + \frac{\mu}{2} \|x\|^2,$$

which is based on the elastic-net regularization [63] leading to nonsmooth strongly convex problems.

For each formulation, we consider a training set $(a_i, b_i)_{i=1}^n$ of n data points, where the b_i 's are scalars in $\{-1, +1\}$ and the a_i 's are feature vectors in \mathbb{R}^d . Then, the goal is to fit a linear model x in \mathbb{R}^d such that the scalar b_i can be predicted well by the inner product $a_i^T x$ or by its sign. Since we normalize the feature vectors a_i , a natural upper bound on the Lipschitz constant L of the unregularized objective can easily be obtained with $L_{\text{logistic}} = 1/4$, $L_{\text{elastic-net}} = 1$, and $L_{\text{lasso}} = 1$.

In the experiments, we consider relatively ill-conditioned problems with the regularization parameter $\mu = 1/(100n)$. The ℓ_1 -regularization parameter is set to $\lambda = 1/n$ for the elastic-net formulation; for the Lasso problem, we consider a logarithmic grid

$10^i/n$, with $i = -3, -2, \dots, 3$, and we select the parameter λ that provides a sparse optimal solution closest to 10% nonzero coefficients.

Data sets. We consider five standard machine learning data sets with different characteristics in terms of size and dimension, which are described in the following table.

Name	covtype	alpha	real-sim	MNIST-CKN	CIFAR-CKN
n	581 012	250 000	72 309	60 000	50 000
d	54	500	20 958	2 304	9 216

The first three data sets are standard machine learning data sets from LIBSVM [13]. We normalize the features to provide a natural estimate of the Lipschitz constant, as mentioned previously. The last two data sets are from computer vision applications. MNIST and CIFAR-10 are two image classification data sets involving 10 classes. The feature representation of each image is computed using an unsupervised convolutional kernel network [39]. We focus here on the task of classifying class #1 vs. other classes.

5.2. Choice of hyper-parameters and variants. We now discuss the choice of default parameters used in the experiments as well as the different variants. First, to deal with the high-dimensional nature of the data, we systematically use the L-BFGS metric H_k and maintain the positive definiteness by skipping updates when necessary (see [20]).

Choice of method \mathcal{M} . We apply QNing to the proximal SVRG algorithm [60] and proximal gradient algorithm. The proximal SVRG algorithm is an incremental algorithm that is able to exploit the finite-sum structure of the objective and can deal with the composite regularization. We also consider the GD algorithm and its proximal variant ISTA, which allows us to perform a comparison with the natural baselines FISTA [2] and L-BFGS.

Stopping criterion for the inner loop. The default stopping criterion consists of solving each subproblem with accuracy $\varepsilon_k \leq \frac{1}{36} \|g_k\|^2$. Although we have shown that such accuracy is attainable within some constant number of iterations, $T = \tilde{O}(n)$, for SVRG with the choice $\kappa = L/2n$, a natural heuristic proposed in Catalyst [34] consists of performing exactly one pass over the data $T = n$ in the inner loop without checking any stopping criterion. In particular, for GD or ISTA, one pass over the data means a single gradient step, because evaluation of the full gradient requires passing through the entire data set. When applying QNing to SVRG and ISTA, we call the default algorithm using the stopping criterion (24) QNing-SVRG and QNing-ISTA, and the one-pass variant QNing-SVRG1 and QNing-ISTA1, respectively.

Choice of regularization parameter κ . We choose κ as in the Catalyst algorithm [34], which is $\kappa = L$ for GD/ISTA and $\kappa = L/2n$ for SVRG. Indeed, the convergence of L-BFGS is hard to characterize, and its theoretical rate of convergence can be pessimistic, as shown in our theoretical analysis. Noting that for smooth functions L-BFGS often outperforms Nesterov's accelerated gradient method, it is reasonable to expect that QNing achieves a similar complexity bound to Catalyst. Later, in section 5.5, we make a comparison between different values of κ to demonstrate the effectiveness of this strategy.

Choice of limited-memory parameter l . The default setting is $l = 100$. Later, in section 5.5, we will compare different values to study the influence of this parameter.

Implementation of the line search. As mentioned earlier, we consider the step sizes η_k in the set $\{1, 1/2, 1/4, 1/8, 0\}$ and select the largest one that satisfies the descent condition.

Evaluation metric. For all experiments, we use the number of gradient evaluations as a measure of complexity, assuming this is the computational bottleneck of all methods considered. This is indeed the case here since the L-BFGS step costs $O(ld)$ floating-point operations [50], whereas evaluating the gradient of the full objective costs $O(nd)$, with $l \ll n$.

5.3. QNing-SVRG for minimizing large sums of functions. We now apply QNing to SVRG and compare different variants.

- SVRG: the Prox-SVRG algorithm of [60] with default parameters $m = 1$ and $\eta = 1/L$, where L is the upper bound on Lipschitz constant of the gradient, as described in section 5.1.
- Catalyst-SVRG: the Catalyst meta-algorithm of [34] applied to Prox-SVRG, using the variant (C3) that performs best among the different variants of Catalyst.
- L-BFGS/Orthant: since implementing L-BFGS effectively with a line-search algorithm is a bit involved, we use the implementation by Mark Schmidt,² which has been widely used in other comparisons [56]. In particular, the Orthant-wise method follows the algorithm developed in [1]. We use L-BFGS for the logistic regression experiment and the Orthant-wise method [1] for elastic-net and Lasso experiments. The limited-memory parameter l is set to 100.
- QNing-SVRG: the algorithm according to the theory given by solving the subproblems until $\varepsilon_k \leq \frac{1}{36} \|g_k\|^2$.
- QNing-SVRG1: the one-pass heuristic.

The result of the comparison is presented in Figure 1 and leads to the conclusions below,³ showing that QNing-SVRG1 is a safe heuristic, which never decreases the speed of the SVRG method.

- L-BFGS/Orthant is less competitive than other approaches that exploit the sum structure of the objective, except on the data set `real-sim`; the difference in performance with the SVRG-based approaches can be important (see data set `alpha`).
- QNing-SVRG1 is significantly faster than or on par with both SVRG and QNing-SVRG.
- QNing-SVRG is significantly faster than, on par with, or only slightly slower than SVRG.
- QNing-SVRG1 is significantly faster than, or on par with Catalyst-SVRG. This justifies our choice of κ , which assumes a priori that L-BFGS performs as well as Nesterov's method.

So far, we have shown that applying QNing with SVRG provides a significant speedup compared to the original SVRG algorithm or other acceleration scheme such as Catalyst. Now we compare our algorithm to other variable metric approaches, including proximal L-BFGS [31] and stochastic L-BFGS [44].

- Proximal L-BFGS: we apply the MATLAB package PNOPT⁴ implemented by [31]. The subproblems are solved by the default algorithm up to the desired accuracy. We consider one subproblem as one gradient evaluation in our plot, even though it often requires multiple passes.

²Available at <http://www.cs.ubc.ca/~schmidtm/Software/minFunc.html>.

³Color figures are available in the online version of this paper.

⁴Available at <https://web.stanford.edu/group/SOL/software/pnopt>.

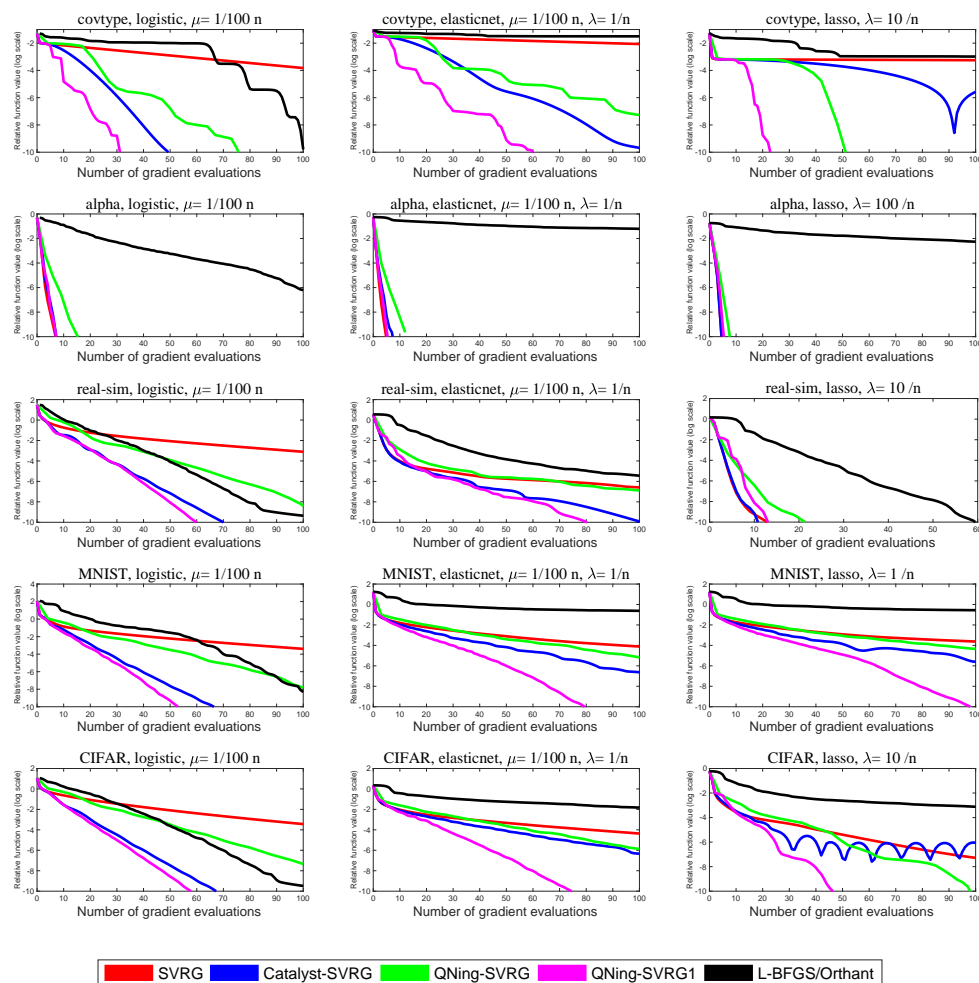


FIG. 1. *Experimental study of the performance of QNing-SVRG for minimizing large sums of functions. We plot the value $F(x_k)/F^* - 1$ as a function of the number of gradient evaluations, on a logarithmic scale; the optimal value F^* is estimated with a duality gap.*

- Stochastic L-BFGS (for smooth objectives): we apply the MATLAB package StochBFGS⁵ implemented by [44]. We consider the “prev” variant, which has the best practical performance.

The result of the comparison is presented in Figure 2 and we observe that QNing-SVRG1 is significantly faster than proximal L-BFGS and stochastic L-BFGS:

- proximal L-BFGS often outperforms Orthant-based methods but it is less competitive than QNing;
- stochastic L-BFGS is sensitive to parameters and data since the variable metric is based on stochastic information that may have high variance. It performs well on data set covtype but becomes less competitive on other data sets. Moreover, it only applies to smooth problems.

⁵Available at <https://perso.telecom-paristech.fr/rgower/software.html>.

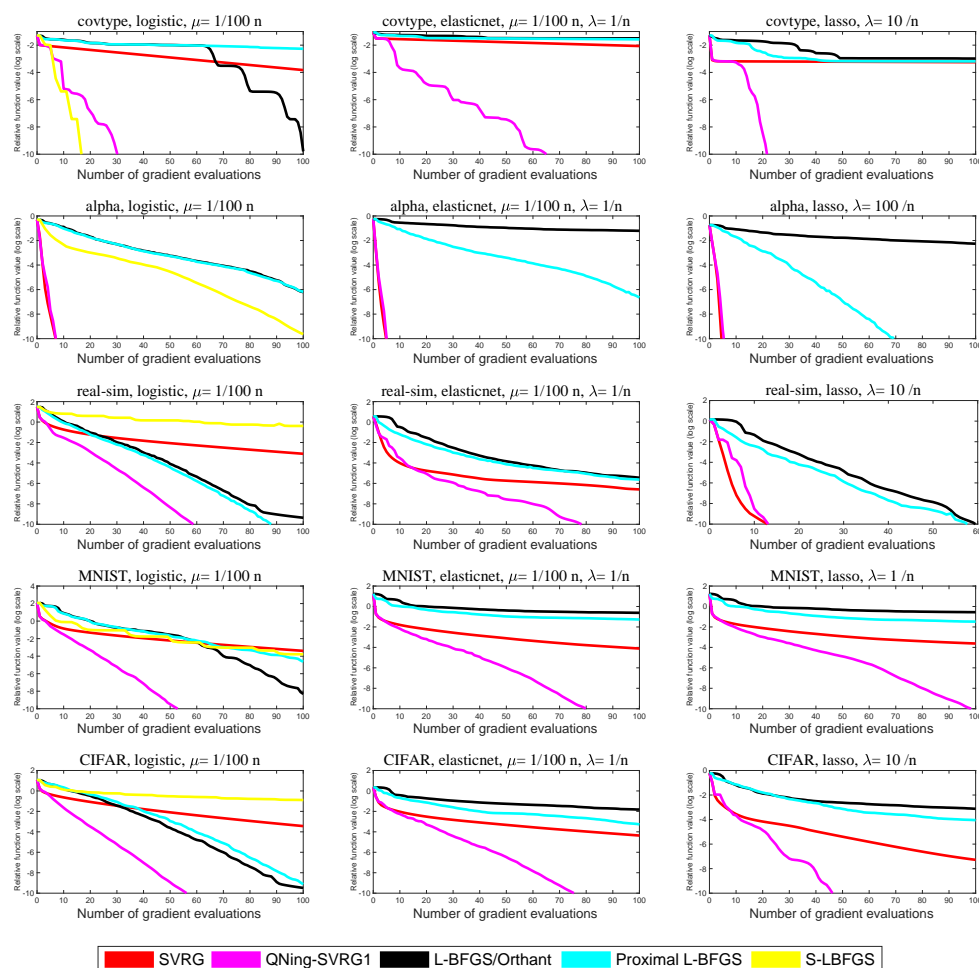


FIG. 2. Comparison to proximal L-BFGS and stochastic L-BFGS. We plot the value $F(x_k)/F^* - 1$ as a function of the number of gradient evaluations, on a logarithmic scale; the optimal value F^* is estimated with a duality gap.

The previous results are complemented by Appendix C.1, which also presents some comparisons in terms of outer-loop iterations, regardless of the cost of the inner loop.

5.4. QNing-ISTA and comparison with L-BFGS. The previous experiments included a comparison between L-BFGS and approaches that are able to exploit the sum structure of the objective. It is interesting to next study the behavior of QNing when applied to a basic proximal gradient descent algorithm such as ISTA. Specifically, we now consider the following:

- GD/ISTA, the classical proximal gradient descent algorithm ISTA [2] with back-tracking line search to automatically adjust the Lipschitz constant of the gradient objective;
- Acc-GD/FISTA, the accelerated variant of ISTA from [2];
- QNing-ISTA and QNing-ISTA1, as in the previous section, replacing SVRG by GD/ISTA.

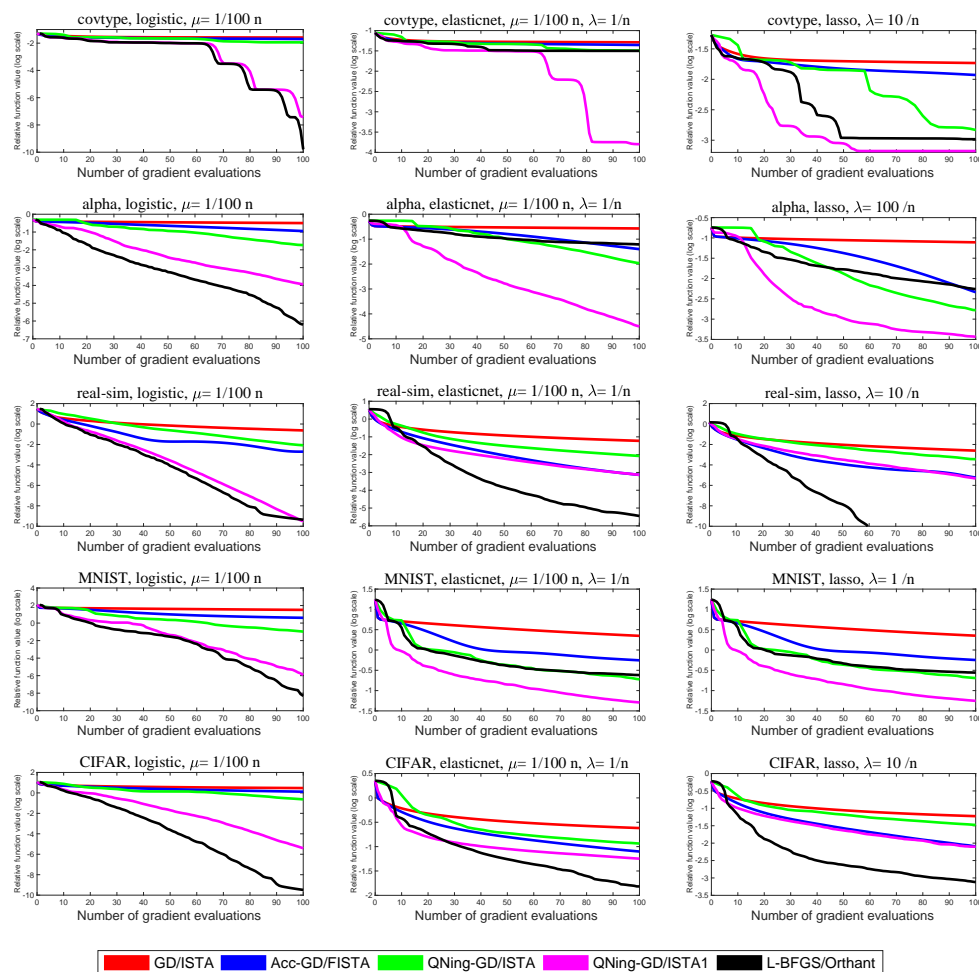


FIG. 3. *Experimental study of the performance of QNing-ISTA. We plot the value $F(x_k)/F^* - 1$ as a function of the number of gradient evaluations, on a logarithmic scale; the optimal value F^* is estimated with a duality gap.*

The results are reported in Figure 3 and lead to the following conclusions.

- L-BFGS is slightly better on average than QNing-ISTA1 for smooth problems, which is not surprising since we use a state-of-the-art implementation with a well-calibrated line search.
- QNing-ISTA1 is always significantly faster than ISTA and QNing-ISTA.
- The QNing-ISTA approaches are significantly faster than FISTA in 12 cases out of 15.
- There is no clear conclusion regarding the performance of the Orthant-wise method versus other approaches. For three data sets—covtype, alpha, and mnist—QNing-ISTA is significantly better than Orthant-wise. However, on the other two data sets, the behavior is different: the Orthant-wise method outperforms QNing-ISTA.

5.5. Experimental study of hyper-parameters l and κ . In this section, we study the influence of the limited-memory parameter l and of the regularization

parameter κ in QNing. More precisely, we start with the parameter l and try the QNing-SVRG1 method with the values $l = 1, 2, 5, 10, 20, 100$. Note that all previous experiments were conducted with $l = 100$, which is the most expensive in terms of memory and computational cost for the L-BFGS step. The results are presented in Figure 4. Interestingly, the experiment suggests that having a large value for l is not necessarily the best choice, especially for composite problems where the solution is sparse, where $l = 10$ seems to perform reasonably well.

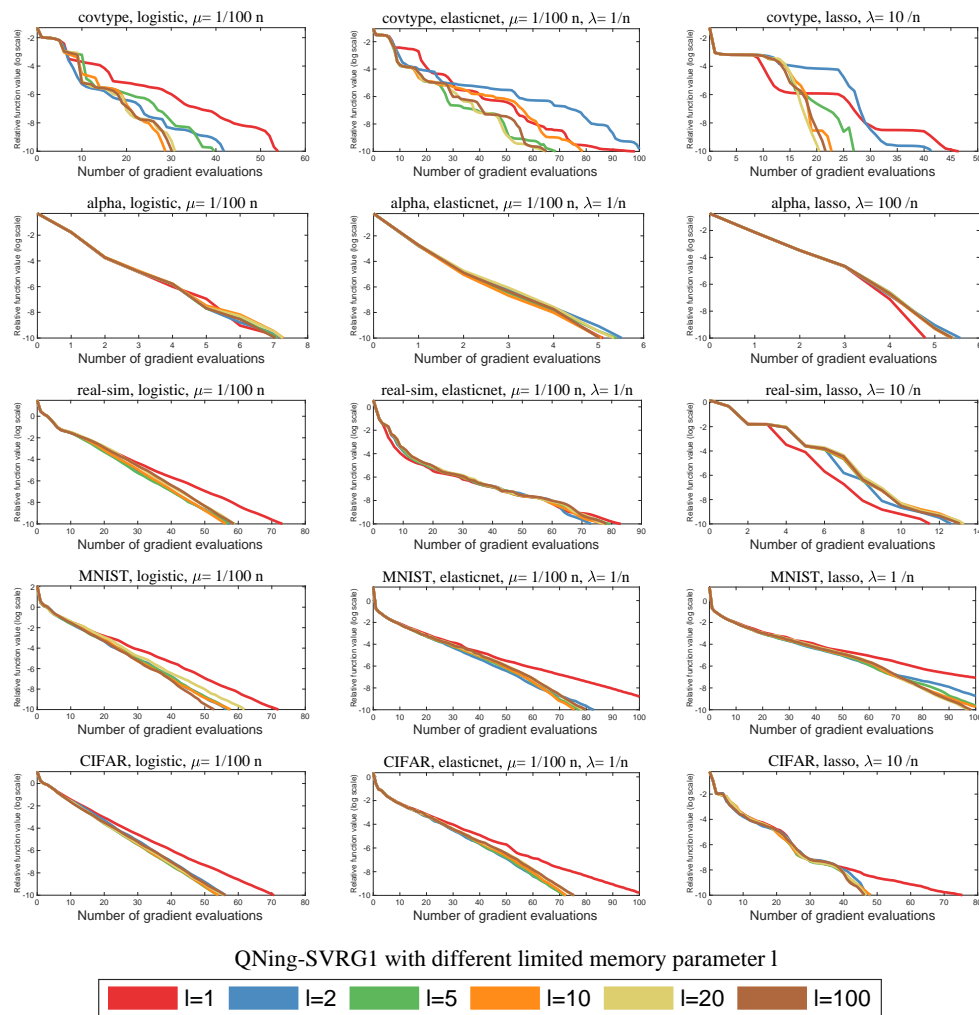


FIG. 4. *Experimental study of influence of the limited-memory parameter l for QNing-SVRG1. We plot the value $F(x_k)/F^* - 1$ as a function of the number of gradient evaluations, on a logarithmic scale; the optimal value F^* is estimated with a duality gap.*

The next experiment consists of studying the robustness of QNing to the smoothing parameter κ . We present in Figure 5 an experiment trying the values $\kappa = 10^i \kappa_0$ for $i = -3, -2, \dots, 2, 3$, where κ_0 is the default parameter used in the previous experiments. The conclusion is clear: QNing clearly slows down when using a larger smoothing parameter than κ_0 , but it is robust to small values of κ (and in fact it even performs better for smaller values than κ_0).

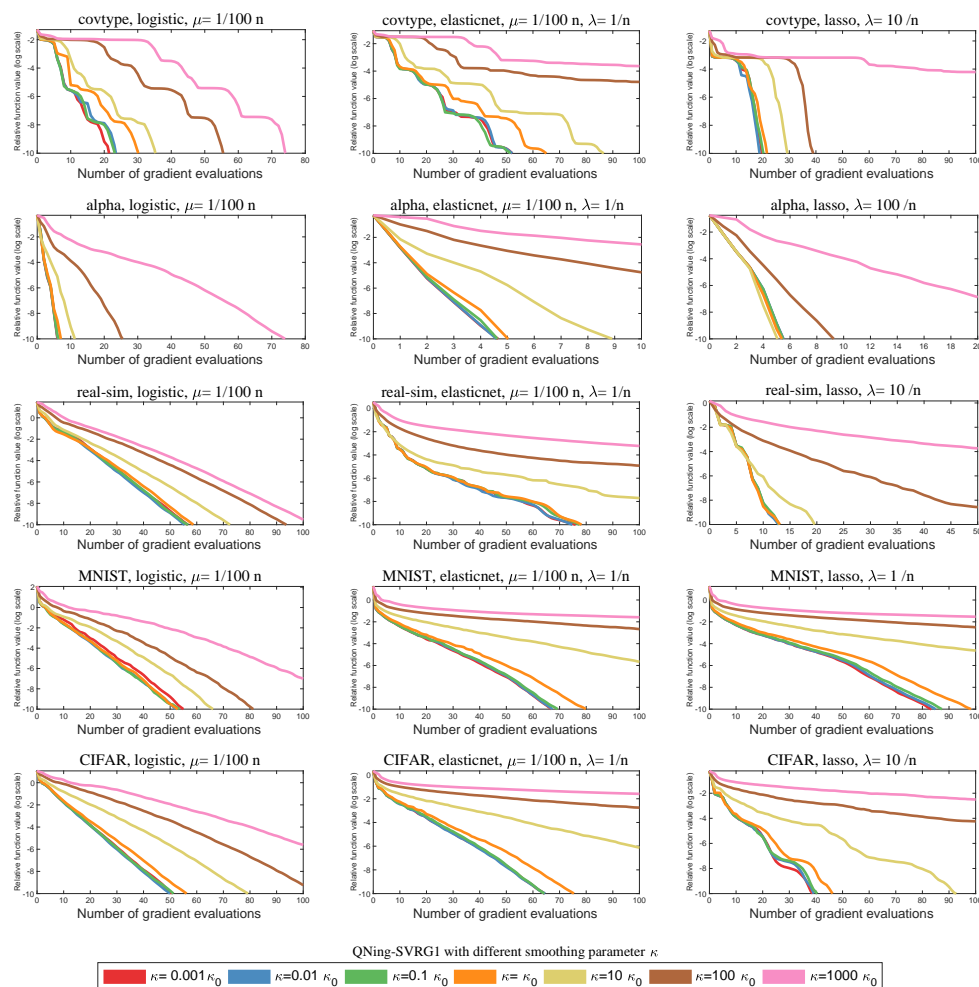


FIG. 5. Experimental study of influence of the smoothing parameter κ for QNing-SVRG1. κ_0 denotes the default choice used in the previous experiments. We plot the value $F(x_k)/F^* - 1$ as a function of the number of gradient evaluations, on a logarithmic scale; the optimal value F^* is estimated with a duality gap.

6. Discussions and concluding remarks. A few questions naturally arise regarding the QNing scheme: one may wonder whether or not our convergence rates may be improved, or if the Moreau envelope could be replaced by another smoothing technique. In this section, we discuss these two points and present concluding remarks.

6.1. Discussion of convergence rates. In this paper, we have established the linear convergence of QNing for strongly convex objectives when subproblems are solved with enough accuracy. Since QNing uses quasi-Newton steps, one might have expected a superlinear convergence rate as several quasi-Newton algorithms often enjoy [11]. The situation is as follows. Consider the BFGS algorithm (without limited memory), as shown in [14]; if the sequence $(\varepsilon_k)_{k \geq 0}$ decreases superlinearly, then it is possible to design a scheme similar to QNing that indeed enjoys a superlinear convergence rate. There is a major downside though: a superlinearly decreasing sequence $(\varepsilon_k)_{k \geq 0}$ implies an exponentially growing number of iterations in the inner

loops, which will degrade the global complexity of the algorithm. This issue makes the approach proposed in [14] impractical.

Another potential strategy for obtaining a faster convergence rate consists in interleaving a Nesterov-type extrapolation step in the QNing algorithm. Indeed, the convergence rate of QNing scales linearly in the condition number μ_F/L_F , which suggests that a faster convergence rate could be obtained using a Nesterov-type acceleration scheme. Empirically, we did not observe any benefit of such a strategy, probably because of the pessimistic nature of the convergence rates that are typically obtained for quasi-Newton approaches based on L-BFGS. Obtaining a linear convergence rate for an L-BFGS algorithm is still an important sanity check, but to the best of our knowledge the gap in performance between these worst-case rates and practice has always been huge for this class of algorithms.

6.2. Other types of smoothing. The Moreau envelope we considered is a particular instance of infimal convolution smoothing [3], whose family also encompasses the so-called Nesterov smoothing [3]. Other ways to smooth a function include randomization techniques [18] or specific strategies tailored for the objective at hand.

One of the main purposes of applying the Moreau envelope is to provide a better conditioning. As recalled in Proposition 1, the gradient of the smoothed function F is κ -Lipschitz continuous regardless of whether the original function is continuously differentiable or not. Furthermore, the conditioning of F is improved with respect to the original function, with a condition number depending on the amount of smoothing. As highlighted in [3], this property is also shared by other types of infimal convolutions. Therefore, QNing could potentially be extended to such types of smoothing in place of the Moreau envelope. A major advantage of our approach, though, is its outstanding simplicity.

6.3. Concluding remarks. To conclude, we have proposed a generic mechanism, QNing, to accelerate existing first-order optimization algorithms with quasi-Newton-type methods. QNing's main features are compatibility with the variable metric update rule and composite optimization. Its ability to combine with incremental approaches makes it a promising tool for solving large-scale machine learning problems. However, a few questions remain open regarding the use of the method in a purely stochastic optimization setting, and the gap in performance between worst-case convergence analysis and practice is significant. We are planning to address the first question about stochastic optimization in future work; the second question is unfortunately difficult and is probably one of the main open questions in the literature about L-BFGS methods.

Appendix A. Proof of Proposition 3. First, we show that the Moreau envelope F inherits the bounded level set property from f .

DEFINITION 2. *We say that a convex function f has bounded level sets if f attains its minimum at x^* in \mathbb{R}^d and, for any x , there exists $R_x > 0$ such that*

$$\forall y \in \mathbb{R}^d \quad \text{s.t.} \quad f(y) \leq f(x) \quad \text{we have} \quad \|y - x^*\| \leq R_x.$$

LEMMA 11. *If f has bounded level sets, then its Moreau envelope F has bounded level sets as well.*

Proof. First, from Proposition 1, the minimum of F is attained at x^* . Next, we reformulate the bounded level set property by contraposition: for any x , there

exists $R_x > 0$ such that

$$\forall y \in \mathbb{R}^d \quad \text{s.t.} \quad \|y - x^*\| > R_x \quad \text{we have} \quad f(y) > f(x).$$

Given x in \mathbb{R}^d , we show that

$$\forall y \in \mathbb{R}^d \quad \text{s.t.} \quad \|y - x^*\| > \sqrt{\frac{2(f(x) - f^*)}{\kappa}} + R_x \quad \text{we have} \quad F(y) > F(x).$$

Let y satisfy the above inequality. By definition,

$$F(y) = f(p(y)) + \frac{\kappa}{2} \|p(y) - y\|^2.$$

From the triangle inequality,

$$\|y - p(y)\| + \|p(y) - x^*\| \geq \|y - x^*\| > \sqrt{\frac{2(f(x) - f^*)}{\kappa}} + R_x.$$

Then either $\|y - p(y)\| > \sqrt{\frac{2(f(x) - f^*)}{\kappa}}$ or $\|p(y) - x^*\| > R_x$.

- If $\|y - p(y)\| > \sqrt{\frac{2(f(x) - f^*)}{\kappa}}$, then

$$F(y) = f(p(y)) + \frac{\kappa}{2} \|p(y) - y\|^2 > f(p(y)) + f(x) - f^* \geq f(x) \geq F(x).$$

- If $\|p(y) - x^*\| > R_x$, then

$$F(y) = f(p(y)) + \frac{\kappa}{2} \|p(y) - y\|^2 \geq f(p(y)) > f(x) \geq F(x).$$

This completes the proof. \square

We are now in a position to prove the proposition.

Proof of Proposition 3. From (25), we have

$$F(x_{k+1}) \leq F(x_k) - \frac{1}{32\kappa} \|\nabla F(x_k)\|^2.$$

Thus, $F(x_k)$ is decreasing. From the bounded level set property of F , there exists $R > 0$ such that $\|x_k - x^*\| \leq R$ for any k . By the convexity of F , we have

$$F(x_k) - F^* \leq \langle \nabla F(x_k), x_k - x^* \rangle \leq \|\nabla F(x_k)\| \|x_k - x^*\| \leq R \|\nabla F(x_k)\|.$$

Therefore,

$$\begin{aligned} F(x_{k+1}) - F^* &\leq F(x_k) - F^* - \frac{1}{32\kappa} \|\nabla F(x_k)\|^2 \\ &\leq F(x_k) - F^* - \frac{(F(x_k) - F^*)^2}{32\kappa R^2}. \end{aligned}$$

Let us define $r_k \triangleq F(x_k) - F^*$. Thus,

$$\frac{1}{r_{k+1}} \geq \frac{1}{r_k(1 - \frac{r_k}{32\kappa R^2})} \geq \frac{1}{r_k} \left(1 + \frac{r_k}{32\kappa R^2}\right) = \frac{1}{r_k} + \frac{1}{32\kappa R^2}.$$

Then, after exploiting the telescoping sum, we obtain

$$\frac{1}{r_{k+1}} \geq \frac{1}{r_0} + \frac{k+1}{32\kappa R^2} \geq \frac{k+1}{32\kappa R^2}. \quad \square$$

Appendix B. Proof of Lemma 8. Let us define $\delta_k = -B_k^{-1}g_k$ and let the subproblems be solved to accuracy $\varepsilon_k \leq \frac{c}{\kappa}\|g_k\|^2$. We show that when $c \leq \frac{\mu^2}{128(\mu+\kappa)^2}$ the following two inequalities hold:

$$(42) \quad F(x_k + \delta_k) \leq F(x_k) - \frac{3}{8\kappa}\|g_k\|^2 + o(\|g_k\|^2)$$

and

$$(43) \quad F_{k+1} \leq F(x_k + \delta_k) + \frac{1}{16\kappa}\|g_k\|^2 + o(\|g_k\|^2).$$

Then, summing the above inequalities yields

$$\begin{aligned} F_{k+1} &\leq F(x_k) - \frac{5}{16\kappa}\|g_k\|^2 + o(\|g_k\|^2) \\ &\leq F_k - \frac{1}{4\kappa}\|g_k\|^2, \end{aligned}$$

where the last inequality holds since $F(x_k) \leq F_k$ and $o(\|g_k\|^2) \leq \frac{1}{4\kappa}\|g_k\|^2$ when k is large enough. This is the desired descent condition (11).

We first prove (42), which relies on the smoothness and Lipschitz Hessian assumption of F . More concretely,

$$\begin{aligned} &F(x_k + \delta_k) - F(x_k) \\ &\leq \nabla F(x_k)^T \delta_k + \frac{1}{2} \delta_k^T \nabla^2 F(x_k) \delta_k + \frac{L_2}{6} \|\delta_k\|^3 \\ &= (\nabla F(x_k) - g_k)^T \delta_k + g_k^T \delta_k + \frac{1}{2} \delta_k^T (\nabla^2 F(x_k) - B_k) \delta_k + \underbrace{\frac{1}{2} \delta_k^T B_k \delta_k}_{= -\frac{1}{2} g_k^T \delta_k} + \frac{L_2}{6} \|\delta_k\|^3 \\ &= \underbrace{\frac{1}{2} g_k^T \delta_k}_{E_1} + \underbrace{(\nabla F(x_k) - g_k)^T \delta_k}_{E_2} + \underbrace{\frac{1}{2} \delta_k^T (\nabla^2 F(x_k) - B_k) \delta_k}_{E_3} + \underbrace{\frac{L_2}{6} \|\delta_k\|^3}_{E_4}. \end{aligned}$$

We shall upper bound each term one by one. First,

$$\begin{aligned} E_1 &= \frac{1}{2} g_k^T \delta_k = -\frac{1}{2} g_k B_k^{-1} g_k \\ &= -\frac{1}{2} g_k \nabla^2 F(x^*)^{-1} g_k - \frac{1}{2} g_k (B_k^{-1} - \nabla^2 F(x^*)^{-1}) g_k \\ &\leq -\frac{1}{2\kappa} \|g_k\|^2 + o(\|g_k\|^2), \end{aligned}$$

where the last inequality uses (DM) and the κ -smoothness of F , which implies

$$\nabla^2 F(x^*) \preceq \kappa I.$$

Second,

$$\begin{aligned} E_2 &= (\nabla F(x_k) - g_k)^T \delta_k \leq \|\nabla F(x_k) - g_k\| \|\delta_k\| \\ &\leq \sqrt{2c} \|g_k\| \|B_k^{-1} g_k\| \quad (\text{from (18)}) \\ &\leq \sqrt{2c} \|g_k\| [\|\nabla^2 F(x^*)^{-1} g_k\| + \|(B_k^{-1} - \nabla^2 F(x^*)^{-1}) g_k\|] \\ &\leq \sqrt{2c} \frac{1}{\mu_F} \|g_k\|^2 + o(\|g_k\|^2) \\ &= \frac{1}{8\kappa} \|g_k\|^2 + o(\|g_k\|^2). \end{aligned}$$

Third,

$$\begin{aligned}
 E_3 &= \frac{1}{2} \delta_k^T (\nabla^2 F(x_k) - B_k) \delta_k \leq \frac{1}{2} \|\delta_k\| \|(\nabla^2 F(x_k) - B_k) \delta_k\| \\
 &\leq \frac{1}{2} \|\delta_k\| (\|(\nabla^2 F(x_k) - \nabla^2 F(x^*)) \delta_k\| + \|(\nabla^2 F(x^*) - B_k) \delta_k\|) \\
 &\leq \frac{L_2}{2} \|x_k - x^*\| \|\delta_k\|^2 + \|\nabla^2 F(x^*)\| \|(B_k^{-1} - \nabla^2 F(x^*)^{-1}) g_k\| \\
 &= o(\|g_k\|^2),
 \end{aligned}$$

where the last line comes from (DM) and the fact that $\|x_k - x^*\| \rightarrow 0$. Last, since

$$\delta_k = \underbrace{-\nabla^2 F(x^*)^{-1} g_k}_{=O(\|g_k\|)} + \underbrace{(\nabla^2 F(x^*)^{-1} - B_k^{-1}) g_k}_{=o(\|g_k\|) \text{ by the Dennis-Moré condition}}$$

and $\|g_k\| \rightarrow 0$, we have

$$E_4 = \frac{L_2}{6} \|\delta_k\|^3 = o(\|g_k\|^2).$$

Summing the above four inequalities yields (42). Next, we prove the other desired inequality, (43). The main effort is to bound $\|g_{k+1}\|$ by a constant factor times $\|g_k\|$. From the inexactness of the subproblem, we have

$$F_{k+1} \leq F(x_{k+1}) + \frac{c}{\kappa} \|g_{k+1}\|^2 \leq F(x_{k+1}) + \frac{2c}{(1-4c)\kappa} \|\nabla F(x_{k+1})\|^2.$$

Moreover,

$$\begin{aligned}
 &\nabla F(x_{k+1}) - \nabla F(x_k) - \nabla^2 F(x^*)(x_{k+1} - x_k) \\
 &= \int_0^1 (\nabla^2 F(x_k + \tau(x_{k+1} - x_k)) - \nabla^2 F(x^*)) (x_{k+1} - x_k) d\tau.
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 &\|\nabla F(x_{k+1}) - \nabla F(x_k) - \nabla^2 F(x^*)(x_{k+1} - x_k)\| \\
 &\leq \max\{\|x_k - x^*\|, \|x_{k+1} - x^*\|\} \|x_{k+1} - x_k\| = o(\|g_k\|).
 \end{aligned}$$

Then, by the triangle inequality, we have

$$\begin{aligned}
 \|\nabla F(x_{k+1})\| &\leq \|\nabla F(x_k) + \nabla^2 F(x^*)(x_{k+1} - x_k)\| + o(\|g_k\|) \\
 &\leq \|\nabla F(x_k) - g_k\| + \underbrace{\|g_k + \nabla^2 F(x^*)(x_{k+1} - x_k)\|}_{=o(\|g_k\|) \text{ by the Dennis-Moré condition}} + o(\|g_k\|) \\
 &\leq \sqrt{2c} \|g_k\| + o(\|g_k\|).
 \end{aligned}$$

As a result,

$$\begin{aligned}
 F_{k+1} &\leq F(x_{k+1}) + \frac{4c^2}{(1-4c)\kappa} \|\nabla g_k\|^2 + o(\|g_k\|^2) \\
 &\leq F(x_{k+1}) + \frac{1}{16\kappa} \|\nabla g_k\|^2 + o(\|g_k\|^2) \quad \text{when } c \leq \frac{1}{16}.
 \end{aligned}$$

This completes the proof. \square

Appendix C. Additional experiments. In this section, we provide additional experimental results, including experimental comparisons in terms of outer loop iterations and an empirical study regarding the choice of the unit step size $\eta_k = 1$.

C.1. Comparisons in terms of outer-loop iterations. In the main paper, we have used the number of gradient evaluations as a natural measure of complexity. Here, we also provide a comparison in terms of outer-loop iterations, which does not take into account the complexity of solving the subproblems. While interesting, the comparison artificially gives an advantage to the stopping criterion (15), since achieving it usually requires multiple passes.

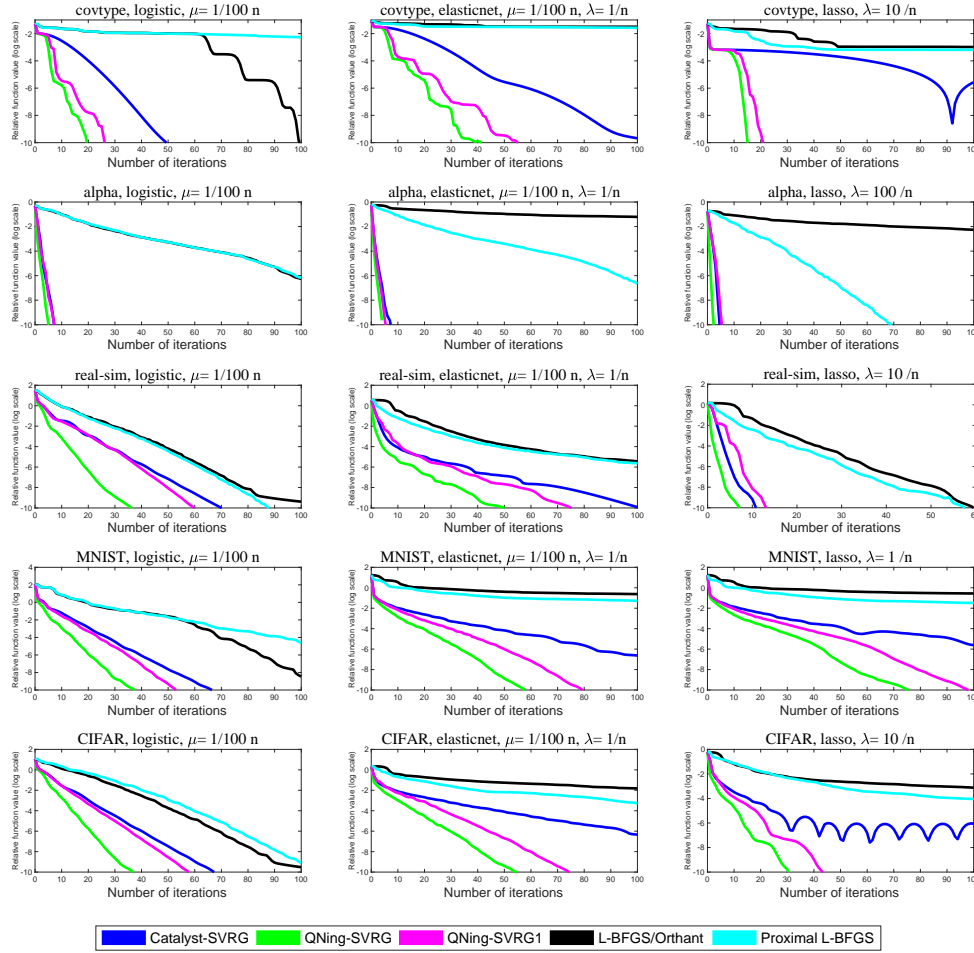


FIG. 6. Experimental study of the performance of QNing-SVRG with respect to the number of outer iterations.

The result of the comparison is presented in Figure 6. We observe that the theory-based variant QNing-SVRG always outperforms the one-pass heuristic QNing-SVRG1. This is not surprising since the subproblems are solved more accurately in the theoretical grounded variant. However, once we take the complexity of the subproblems into account, QNing-SVRG never outperforms QNing-SVRG1. This suggests that it is not beneficial to solve the subproblem up to high accuracy as long as the algorithms converge.

C.2. Empirical frequency of choosing the unit step size. In this section, we evaluate how often the unit step size is taken in the line search. When the unit step size is taken, the variable metric step provides a sufficient decrease, which is key for acceleration. The statistics of QNing-SVRG1 (the one-pass variant) and QNing-SVRG (the subproblems are solved until the stopping criterion (15) is satisfied) are given in Tables 1 and 2, respectively. As we can see, for most of the iterations ($> 90\%$), the unit step size is taken.

TABLE 1

Relative frequency of picking the unit step size $\eta_k = 1$ of QNing-SVRG1. For each method, the first column is in the form N/D , where N is the number of times the unit step size was picked over the iterations and D is the total number of iterations. The total number of iterations D varies a lot since we stop our algorithm as soon as the relative function gap is smaller than 10^{-10} or the maximum number of iterations (100) is reached. It implicitly indicates how easy the problem is.

	Logistic		Elastic-net		Lasso	
covtype	24/27	89%	54/56	96%	19/21	90%
alpha	8/8	100%	6/6	100%	6/6	100%
real-sim	60/60	100%	71/76	93%	14/14	100%
MNIST	53/53	100%	80/80	100%	100/100	100%
CIFAR-10	58/58	100%	75/75	100%	42/44	95%

TABLE 2

Relative frequency of picking the unit step size $\eta_k = 1$ of QNing-SVRG. The settings are the same as in Table 1.

	Logistic		Elastic-net		Lasso	
covtype	18/20	90%	23/25	92%	16/16	100%
alpha	6/6	100%	4/4	100%	3/3	100%
real-sim	27/27	100%	20/23	87%	8/8	100%
MNIST	27/27	100%	28/28	100%	28/28	100%
CIFAR-10	25/25	100%	29/29	100%	31/31	100%

Acknowledgment. The authors would like to thank the editor and the reviewers for their constructive and detailed comments.

REFERENCES

- [1] G. ANDREW AND J. GAO, *Scalable training of L_1 -regularized log-linear models*, in Proceedings of the 24th International Conference on Machine Learning, Association for Computing Machinery, New York, 2007, pp. 33–40.
- [2] A. BECK AND M. TEOULLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*, SIAM J. Imaging Sci., 2 (2009), pp. 183–202.
- [3] A. BECK AND M. TEOULLE, *Smoothing and first order methods: A unified framework*, SIAM J. Optim., 22 (2012), pp. 557–580.
- [4] S. BECKER AND J. FADILI, *A quasi-Newton proximal splitting method*, in Adv. Neural Inf. Process. Syst. 25, Curran Associates, Red Hook, NY, 2012, pp. 2618–2626.
- [5] D. P. BERTSEKAS, *Convex Optimization Algorithms*, Athena Scientific, Nashua, NH, 2015.
- [6] J.-F. BONNANS, J. C. GILBERT, C. LEMARÉCHAL, AND C. A. SAGASTIZÁBAL, *Numerical Optimization: Theoretical and Practical Aspects*, 2nd edn, Springer, Berlin, 2006.
- [7] J. BURKE AND M. QIAN, *On the superlinear convergence of the variable metric proximal point algorithm using Broyden and BFGS matrix secant updating*, Math. Program., 88 (2000), pp. 157–181.
- [8] R. H. BYRD, G. M. CHIN, J. NOCEDAL, AND Y. WU, *Sample size selection in optimization methods for machine learning*, Math. Program., 134 (2012), pp. 127–155.

- [9] R. H. BYRD, S. L. HANSEN, J. NOCEDAL, AND Y. SINGER, *A stochastic quasi-Newton method for large-scale optimization*, SIAM J. Optim., 26 (2016), pp. 1008–1031.
- [10] R. H. BYRD, J. NOCEDAL, AND F. OZTOPRAK, *An inexact successive quadratic approximation method for $L - 1$ regularized optimization*, Math. Program., 157 (2015), pp. 375–396.
- [11] R. H. BYRD, J. NOCEDAL, AND Y.-X. YUAN, *Global convergence of a case of quasi-Newton methods on convex problems*, SIAM J. Numer. Anal., 24 (1987), pp. 1171–1190.
- [12] C. CARTIS AND K. SCHEINBERG, *Global convergence rate analysis of unconstrained optimization methods based on probabilistic models*, Math. Program., 169 (2018), pp. 337–375.
- [13] C. CHANG AND C. LIN, *LIBSVM: A library for support vector machines*, ACM Trans. Intell. Syst. Tech., 2 (2011), No. 27.
- [14] X. CHEN AND M. FUKUSHIMA, *Proximal quasi-Newton methods for nondifferentiable convex optimization*, Math. Program., 85 (1999), pp. 313–334.
- [15] A. DEFAZIO, F. BACH, AND S. LACOSTE-JULIEN, *SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives*, in Adv. Neural Inf. Process. Syst. 27, Curran Associates, Red Hook, NY, 2014, pp. 1646–1654.
- [16] A. DEFAZIO, J. DOMKE, AND T. S. CAETANO, *Finito: A faster, permutable incremental gradient method for big data problems*, in Proceedings of the International Conference on Machine Learning 2014, Proc. Mach. Learn. Res. 32, 2014, pp. 1125–1133; available at <http://proceedings.mlr.press/v32/>.
- [17] J. E. DENNIS AND J. J. MORÉ, *A characterization of superlinear convergence and its application to quasi-Newton methods*, Math. Comp., 28 (1974), pp. 549–560.
- [18] J. C. DUCHI, P. L. BARTLETT, AND M. J. WAINWRIGHT, *Randomized smoothing for stochastic optimization*, SIAM J. Optim., 22 (2012), pp. 674–701.
- [19] M. ELAD, *Sparse and Redundant Representations*, Springer, New York, 2010.
- [20] M. P. FRIEDLANDER AND M. SCHMIDT, *Hybrid deterministic-stochastic methods for data fitting*, SIAM J. Sci. Comput., 34 (2012), pp. A1380–A1405.
- [21] R. FROSTIG, R. GE, S. M. KAKADE, AND A. SIDFORD, *Un-regularizing: Approximate proximal point and faster stochastic algorithms for empirical risk minimization*, in Proceedings of the 32nd International Conference on Machine Learning, Proc. Mach. Learn. Res. 37, 2015, pp. 2540–2548; available at <http://proceedings.mlr.press/v37/>.
- [22] M. FUENTES, J. MALICK, AND C. LEMARÉCHAL, *Descentwise inexact proximal algorithms for smooth optimization*, Comput. Optim. Appl., 53 (2012), pp. 755–769.
- [23] M. FUKUSHIMA AND L. QI, *A globally and superlinearly convergent algorithm for nonsmooth convex minimization*, SIAM J. Optim., 6 (1996), pp. 1106–1120.
- [24] S. GHADIMI, G. LAN, AND H. ZHANG, *Generalized Uniformly Optimal Methods for Nonlinear Programming*, preprint, <https://arxiv.org/pdf/1508.07384>, 2015.
- [25] H. GHANBARI AND K. SCHEINBERG, *Proximal quasi-Newton methods for regularized convex optimization with linear and accelerated sublinear convergence rates*, Comput. Optim. Appl., 69 (2018), pp. 597–627.
- [26] R. M. GOWER, D. GOLDFARB, AND P. RICHTÁRIK, *Stochastic block BFGS: Squeezing more curvature out of data*, in Proceedings of the International Conference on Machine Learning 2016, Proc. Mach. Learn. Res. 48, 2016, pp. 1869–1878; available at <http://proceedings.mlr.press/v48/>.
- [27] O. GÜLER, *New proximal point algorithms for convex minimization*, SIAM J. Optim., 2 (1992), pp. 649–664.
- [28] J.-B. HIRIART-URRUTY AND C. LEMARÉCHAL, *Convex Analysis and Minimization Algorithms I*, Grundlehren Math. Wiss. 305, Springer, Berlin, 1996.
- [29] J.-B. HIRIART-URRUTY AND C. LEMARÉCHAL, *Convex Analysis and Minimization Algorithms. II*, Grundlehren Math. Wiss. 306, Springer, Berlin, 1996.
- [30] C. LEE AND S. J. WRIGHT, *Inexact Successive Quadratic Approximation for Regularized Optimization*, preprint, <https://arxiv.org/abs/1803.01298>, 2018.
- [31] J. LEE, Y. SUN, AND M. SAUNDERS, *Proximal Newton-type methods for convex optimization*, in Adv. Neural Inf. Process. Syst. 25, Curran Associates, Red Hook, NY, 2012, pp. 827–835.
- [32] C. LEMARÉCHAL AND C. SAGASTIZÁBAL, *Practical aspects of the Moreau–Yosida regularization: Theoretical preliminaries*, SIAM J. Optim., 7 (1997), pp. 367–385.
- [33] H. LIN, J. MAIRAL, AND Z. HARCHAOUI, *A universal catalyst for first-order optimization*, in Adv. Neural Inf. Process. Syst. 28, Curran Associates, Red Hook, NY, 2015, pp. 3384–3392.
- [34] H. LIN, J. MAIRAL, AND Z. HARCHAOUI, *Catalyst acceleration for first-order convex optimization: From theory to practice*, J. Mach. Learn. Res., 18 (2018), pp. 7854–7907.
- [35] D. C. LIU AND J. NOCEDAL, *On the limited memory BFGS method for large scale optimization*, Math. Program., 45 (1989), pp. 503–528.
- [36] X. LIU, C.-J. HSIEH, J. D. LEE, AND Y. SUN, *An Inexact Subsampled Proximal Newton-Type*

- Method for Large-Scale Machine Learning*, preprint, <https://arxiv.org/pdf/1708.08552>, 2017.
- [37] J. MAIRAL, *Sparse Coding for Machine Learning, Image Processing and Computer Vision*, Ph.D. thesis, École normale supérieure, Cachan, 2010.
 - [38] J. MAIRAL, *Incremental majorization-minimization optimization with application to large-scale machine learning*, SIAM J. Optim., 25 (2015), pp. 829–855.
 - [39] J. MAIRAL, *End-to-end kernel learning with supervised convolutional kernel networks*, in Adv. Neural Inf. Process. Syst. 29, Curran Associates, Red Hook, NY, 2016, pp. 1399–1407.
 - [40] J. MAIRAL, F. BACH, AND J. PONCE, *Sparse modeling for image and vision processing*, Found. Trends Comput. Graph. Vision, 8 (2014), pp. 85–283.
 - [41] R. MIFFLIN, *A quasi-second-order proximal bundle algorithm*, Math. Program., 73 (1996), pp. 51–72.
 - [42] A. MOKHTARI AND A. RIBEIRO, *Global convergence of online limited memory BFGS*, J. Mach. Learn. Res., 16 (2015), pp. 3151–3181.
 - [43] J.-J. MOREAU, *Fonctions convexes duales et points proximaux dans un espace Hilbertien*, C. R. Acad. Sci. Paris Sér. A Math., 255 (1962), pp. 2897–2899.
 - [44] P. MORITZ, R. NISHIHARA, AND M. I. JORDAN, *A linearly-convergent stochastic L-BFGS algorithm*, in Proceedings of the 19th International Conference on Artificial Intelligence and Statistics, Proc. Mach. Learn. Res. 51, 2016, pp. 249–258, available at <http://proceedings.mlr.press/v51/>.
 - [45] Y. NESTEROV, *A method of solving a convex programming problem with convergence rate $O(1/k^2)$* , Soviet Math. Dokl., 27 (1983), pp. 372–376.
 - [46] Y. NESTEROV, *Smooth minimization of non-smooth functions*, Math. Program., 103 (2005), pp. 127–152.
 - [47] Y. NESTEROV, *Efficiency of coordinate descent methods on huge-scale optimization problems*, SIAM J. Optim., 22 (2012), pp. 341–362.
 - [48] Y. NESTEROV, *Gradient methods for minimizing composite functions*, Math. Program., 140 (2013), pp. 125–161.
 - [49] J. NOCEDAL, *Updating quasi-Newton matrices with limited storage*, Math. Comp., 35 (1980), pp. 773–782.
 - [50] J. NOCEDAL AND S. WRIGHT, *Numerical Optimization*, Springer Ser. Oper. Res. Financ. Eng., Springer, New York, 2006.
 - [51] M. RAZAVIYAYN, M. HONG, AND Z.-Q. LUO, *A unified convergence analysis of block successive minimization methods for nonsmooth optimization*, SIAM J. Optim., 23 (2013), pp. 1126–1153.
 - [52] P. RICHTÁRIK AND M. TAKÁČ, *Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function*, Math. Program., 144 (2014), pp. 1–38.
 - [53] R. T. ROCKAFELLAR, *Monotone operators and the proximal point algorithm*, SIAM J. Control Optim., 14 (1976), pp. 877–898.
 - [54] K. SCHEINBERG AND X. TANG, *Practical inexact proximal quasi-Newton method with global complexity analysis*, Math. Program., 160 (2016), pp. 495–529.
 - [55] M. SCHMIDT, D. KIM, AND S. SRA, *Projected Newton-Type Methods in Machine Learning*, MIT Press, Cambridge, MA, 2011, pp. 305–330.
 - [56] M. SCHMIDT, N. L. ROUX, AND F. BACH, *Minimizing finite sums with the stochastic average gradient*, Math. Program., 160 (2017), pp. 83–112.
 - [57] N. N. SCHRAUDOLPH, J. YU, AND S. GÜNTER, *A stochastic quasi-Newton method for online convex optimization*, in Proceedings of the Eleventh International Conference on Artificial Intelligence and Statistics, Proc. Mach. Learn. Res. 2, 2007, pp. 436–443; available at <http://proceedings.mlr.press/v2/>.
 - [58] S. SHALEV-SHWARTZ AND T. ZHANG, *Accelerated proximal stochastic dual coordinate ascent for regularized loss minimization*, Math. Program., 155 (2016), pp. 105–145.
 - [59] L. STELLA, A. THEMELIS, AND P. PATRINOS, *Forward-backward quasi-Newton methods for non-smooth optimization problems*, Comput. Optim. Appl., 67 (2017), pp. 443–487.
 - [60] L. XIAO AND T. ZHANG, *A proximal stochastic gradient method with progressive variance reduction*, SIAM J. Optim., 24 (2014), pp. 2057–2075.
 - [61] K. YOSIDA, *Functional Analysis*, Springer, Berlin, 1980.
 - [62] J. YU, S. VISHWANATHAN, S. GÜNTER, AND N. N. SCHRAUDOLPH, *A quasi-Newton approach to non-smooth convex optimization*, in Proceedings of the 25th International Conference on Machine Learning, Association for Computing Machinery, New York, 2008, pp. 1216–1223.
 - [63] H. ZOU AND T. HASTIE, *Regularization and variable selection via the elastic net*, J. R. Stat. Soc. Ser. B. Stat. Methodol., 67 (2005), pp. 301–320.