

# IQN: AN INCREMENTAL QUASI-NEWTON METHOD WITH LOCAL SUPERLINEAR CONVERGENCE RATE\*

ARYAN MOKHTARI<sup>†</sup>, MARK EISEN<sup>†</sup>, AND ALEJANDRO RIBEIRO<sup>†</sup>

**Abstract.** The problem of minimizing an objective that can be written as the sum of a set of  $n$  smooth and strongly convex functions is challenging because the cost of evaluating the function and its derivatives is proportional to the number of elements in the sum. The Incremental Quasi-Newton (IQN) method proposed here belongs to the family of stochastic and incremental methods that have a cost per iteration independent of  $n$ . IQN iterations are a stochastic version of BFGS iterations that use memory to reduce the variance of stochastic approximations. The method is shown to exhibit local superlinear convergence. The convergence properties of IQN bridge a gap between deterministic and stochastic quasi-Newton methods. Deterministic quasi-Newton methods exploit the possibility of approximating the Newton step using objective gradient differences. They are appealing because they have a smaller computational cost per iteration relative to Newton’s method and achieve a superlinear convergence rate under customary regularity assumptions. Stochastic quasi-Newton methods utilize stochastic gradient differences in lieu of actual gradient differences. This makes their computational cost per iteration independent of the number of objective functions  $n$ . However, existing stochastic quasi-Newton methods have sublinear or linear convergence at best. IQN is the first stochastic quasi-Newton method proven to converge superlinearly in a local neighborhood of the optimal solution. IQN differs from state-of-the-art incremental quasi-Newton methods in three aspects: (i) The use of aggregated information of variables, gradients, and quasi-Newton Hessian approximation matrices to reduce the noise of gradient and Hessian approximations. (ii) The approximation of each individual function by its Taylor’s expansion in which the linear and quadratic terms are evaluated with respect to the same iterate. (iii) The use of a cyclic scheme to update the functions in lieu of a random selection routine. We use these fundamental properties of IQN to establish its local superlinear convergence rate. The presented numerical experiments match our theoretical results and justify the advantage of IQN relative to other incremental methods.

**Key words.** Large-scale optimization, stochastic optimization, quasi-Newton methods, incremental methods, superlinear convergence

**AMS subject classifications.** 90C06, 90C25, 90C30, 90C52

**1. Introduction.** We study large scale optimization problems with objective functions expressed as the sum of a set of components which arise often in application domains such as machine learning [3, 2, 30, 8], control [5, 7, 15], and wireless communications [27, 23, 24]. Formally, we consider a variable  $\mathbf{x} \in \mathbb{R}^p$  and a function  $f$  which is defined as the average of  $n$  smooth and strongly convex functions labelled  $f_i : \mathbb{R}^p \rightarrow \mathbb{R}$  for  $i = 1, \dots, n$ . We refer to individual functions  $f_i$  as sample functions and to the total number of functions  $n$  as the sample size. Our goal is to find the optimal argument  $\mathbf{x}^*$  that solves the strongly convex program

$$(1) \quad \mathbf{x}^* := \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}) := \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}).$$

We restrict attention to cases where the component functions  $f_i$  are strongly convex and their gradients are Lipschitz continuous. We further focus in problems where  $n$  is large enough so as to warrant application of stochastic or iterative methods. Our goal is to propose an iterative quasi-Newton method to solve (1) which is shown to exhibit a local superlinear convergence rate. This is achieved while performing local iterations with a cost of order  $\mathcal{O}(p^2)$  independent of the number of samples  $n$ .

\*This work was supported by ONR N00014-12-1-0997.

<sup>†</sup>Department of Electrical and Systems Engineering, University of Pennsylvania, Philadelphia, PA (aryanm@seas.upenn.edu,maeisen@seas.upenn.edu,aribeiro@seas.upenn.edu)

Setting temporarily aside the complications related to the number of component functions, the minimization of  $f$  in (1) can be carried out using iterative descent algorithms. A simple solution is to use gradient descent (GD) which iteratively descends along gradient directions  $\nabla f(\mathbf{x}) = (1/n) \sum_{i=1}^n \nabla f_i(\mathbf{x})$ . GD incurs a per iteration computational cost of order  $\mathcal{O}(np)$  and is known to converge at a linear rate towards  $\mathbf{x}^*$  under the hypotheses we have placed on  $f$ . Whether the linear convergence rate of GD is acceptable depends on the desired accuracy and on the condition number of  $f$  which, when large, can make the convergence constant close to one. As one or both of these properties often limit the applicability of GD, classical alternatives to improve convergence rates have been developed. Newton’s method adapts to the curvature of the objective by computing Hessian inverses and converges at a quadratic rate in a local neighborhood of the optimal argument irrespective of the problem’s condition number. To achieve this quadratic convergence rate, we must evaluate and invert Hessians resulting in a per iteration cost of order  $\mathcal{O}(np^2 + p^3)$ . Quasi-Newton methods build on the idea of approximating the Newton step using first-order information of the objective function and exhibit local superlinear convergence [4, 22, 12]. An important feature of quasi-Newton methods is that they have a per iteration cost of order  $\mathcal{O}(np + p^2)$ , where the term  $\mathcal{O}(np)$  corresponds to the cost of gradient computation and the cost  $\mathcal{O}(p^2)$  indicates the computational complexity of updating the approximate Hessian inverse matrix.

The combination of a local superlinear convergence rate and the smaller computational cost per iteration relative to Newton – a reduction by a factor of  $p$  operations per iteration – make quasi-Newton methods an appealing choice. In the context of optimization problems having the form in (1), quasi-Newton methods also have the advantage that curvature is estimated using gradient evaluations. To see why this is meaningful we must recall that the customary approach to avoid the  $\mathcal{O}(np)$  computational cost of GD iterations is to replace gradients  $\nabla f(\mathbf{x})$  by their stochastic approximations  $\nabla f_i(\mathbf{x})$ , which can be evaluated with a cost of order  $\mathcal{O}(p)$ . One can then think of using stochastic versions of these gradients to develop stochastic quasi-Newton methods with per iterations cost of order  $\mathcal{O}(p + p^2)$ . This idea was demonstrated to be feasible in [29] which introduces a stochastic (online) version of the BFGS quasi-Newton method as well as a stochastic version of its limited memory variant. Although [29] provides numerical experiments illustrating significant improvements in convergence times relative to stochastic (S) GD, theoretical guarantees are not established.

The issue of proving convergence of stochastic quasi-Newton methods is tackled in [18] and [19]. In [18] the authors show that stochastic BFGS may not be convergent because the Hessian approximation matrices can become close to singular. A regularized stochastic BFGS (RES) method is proposed by changing the proximity condition of BFGS to ensure that the eigenvalues of the Hessian inverse approximation are uniformly bounded. Enforcing this property yields a provably convergent algorithm. In [19] the authors show that the limited memory version of stochastic (online) BFGS proposed in [29] is almost surely convergent and has a sublinear convergence rate in expectation. This is achieved without using regularizations. An alternative provably convergent stochastic quasi-Newton method is proposed in [6]. This method differs from those in [29, 18, 19] in that it collects (stochastic) second order information to estimate the objective’s curvature. This is in contrast to estimating curvature using the difference of two consecutive stochastic gradients.

Although the methods in [29, 18, 19, 6] are successful in expanding the application of quasi-Newton methods to stochastic settings, their convergence rate is sublinear.

This is not better than the convergence rate of SGD and, as is also the case in SGD, is a consequence of the stochastic approximation noise which necessitates the use of diminishing stepsizes. The stochastic quasi-Newton methods in [16, 20] resolve this issue by using the variance reduction technique proposed in [13]. The fundamental idea of the work in [13] is to reduce the noise of the stochastic gradient approximation by computing the exact gradient in an outer loop to use it in an inner loop for gradient approximation. The methods in [16, 20], which incorporate the variance reduction scheme presented in [13] into the update of quasi-Newton methods, are successful in achieving a linear convergence rate.

At this point, we must remark on an interesting mismatch. The convergence rate of SGD is sublinear, and the convergence rate of deterministic GD is linear. The use of variance reduction techniques in SGD recovers the linear convergence rate of GD, [13]. On the other hand, the convergence rate of stochastic quasi-Newton methods is sublinear, and the convergence rate of deterministic quasi-Newton methods is superlinear. The use of variance reduction in stochastic quasi-Newton methods achieves linear convergence but does not recover a superlinear rate. Hence, a fundamental question remains unanswered: Is it possible to design an incremental quasi-Newton method that recovers the superlinear convergence rate of deterministic quasi-Newton algorithms? In this paper, we show that the answer to this open problem is positive by proposing an incremental quasi-Newton method (IQN) with a local superlinear convergence rate. This is the first quasi-Newton method to achieve superlinear convergence while having a per iteration cost independent of the number of functions  $n$  – the cost per iteration is of order  $\mathcal{O}(p^2)$ .

There are three major differences between the IQN method and state-of-the-art incremental (stochastic) quasi-Newton methods that lead to the former’s superlinear convergence rate. First, the proposed IQN method uses the aggregated information of variables, gradients, and Hessian approximation matrices to reduce the noise of approximation for both gradients and Hessian approximation matrices. This is different to the variance-reduced stochastic quasi-Newton methods in [16, 20] that attempt to reduce only the noise of gradient approximations. Second, in IQN the index of the updated function is chosen in a cyclic fashion, rather than the random selection scheme used in the incremental methods in [29, 18, 19, 6]. The cyclic routine in IQN allows to bound the error at each iteration as a function of the errors of the last  $n$  iterates, something that is not possible when using a random scheme. To explain the third and most important difference we point out that the form of quasi-Newton updates is the solution of a local second order Taylor approximation of the objective. It is possible to understand stochastic quasi-Newton methods as an analogous approximation of individual sample functions. However, it turns out that the state-of-the-art stochastic quasi-methods evaluate the linear and quadratic terms of the Taylor’s expansion at different points yielding and inconsistent approximation (Remark 7). The IQN method utilizes a consistent Taylor series which yields a more involved update which we nonetheless show can be implemented with the same computational cost. These three properties together lead to an incremental quasi-Newton method with a local superlinear convergence rate.

We start the paper by recapping the BFGS quasi-Newton method and the Dennis-Moré condition which is sufficient and necessary to prove superlinear convergence rate of the BFGS method (Section 2). Then, we present the proposed Incremental Quasi-Newton method (IQN) as an incremental aggregated version of the traditional BFGS method (Section 3). We first explain the difference between the Taylor’s expansion used in IQN and state-of-the-art incremental (stochastic) quasi-Newton methods. Fur-

ther, we explain the mechanism for aggregation of the functions information and the scheme for updating the stored information. Moreover, we present an efficient implementation of the proposed IQN method with a computational complexity of the order  $\mathcal{O}(p^2)$  (Section 3.1). The convergence analysis of the IQN method is then presented (Section 4). We use the classic analysis of quasi-Newton methods to show that in a local neighborhood of the optimal solution the sequence of variables converges to the optimal argument  $\mathbf{x}^*$  linearly after each pass over the set of functions (Lemma 3). We use this result to show that for each component function  $f_i$  the Dennis-Moré condition holds (Proposition 4). However, this condition is not sufficient to prove superlinear convergence of the sequence of errors  $\|\mathbf{x}^t - \mathbf{x}^*\|$ , since it does not guarantee the Dennis-Moré condition for the global objective  $f$ . To overcome this issue we introduce a novel convergence analysis approach which exploits the local linear convergence of IQN to present a more general version of the Dennis-Moré condition for each component function  $f_i$  (Lemma 5). We exploit this result to establish superlinear convergence of the iterates generated by IQN (Theorem 6). In Section 6, we present numerical simulation results, comparing the performance of IQN to that of first-order incremental and stochastic methods. We test the performance on a set of large-scale regression problems and observe strong numerical gain in total computation time relative to existing methods.

**1.1. Notation.** Vectors are written as lowercase  $\mathbf{x} \in \mathbb{R}^p$  and matrices as uppercase  $\mathbf{A} \in \mathbb{R}^{p \times p}$ . We use  $\|\mathbf{x}\|$  and  $\|\mathbf{A}\|$  to denote the Euclidean norm of vector  $\mathbf{x}$  and matrix  $\mathbf{A}$ , respectively. Given a positive definite matrix  $\mathbf{M}$ , the weighted matrix norm  $\|\mathbf{A}\|_{\mathbf{M}}$  is defined as  $\|\mathbf{A}\|_{\mathbf{M}} := \|\mathbf{M}\mathbf{A}\mathbf{M}\|_{\mathbf{F}}$ , where  $\|\cdot\|_{\mathbf{F}}$  is the Frobenius norm. Given a function  $f$  its gradient and Hessian at point  $\mathbf{x}$  are denoted as  $\nabla f(\mathbf{x})$  and  $\nabla^2 f(\mathbf{x})$ , respectively.

**2. BFGS Quasi-Newton Method.** Consider the problem in (1) for relatively large  $n$ . In a conventional optimization setting, this can be solved using a quasi-Newton method that iteratively updates a variable  $\mathbf{x}^t$  for  $t = 0, 1, \dots$  based on the general recursive expression

$$(2) \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \eta^t (\mathbf{B}^t)^{-1} \nabla f(\mathbf{x}^t),$$

where  $\eta^t$  is a scalar stepsize and  $\mathbf{B}^t$  is a positive definite matrix that approximates the exact Hessian of the objective function  $\nabla^2 f(\mathbf{x}^t)$ . The stepsize  $\eta^t$  is evaluated based on a line search routine for the global convergence of quasi-Newton methods. Our focus in this paper, however, is on the local convergence of quasi-Newton methods, which requires the unit stepsize  $\eta^t = 1$ . Therefore, throughout the paper we assume that the variable  $\mathbf{x}^t$  is close to the optimal solution  $\mathbf{x}^*$  – we will formalize the notion of being close to the optimal solution – and the stepsize is  $\eta^t = 1$ .

The goal of quasi-Newton methods is to compute the Hessian approximation matrix  $\mathbf{B}^t$  and its inverse  $(\mathbf{B}^t)^{-1}$  by using only the first-order information, i.e., gradients, of the objective. Their use is widespread due to the many applications in which the Hessian information required in Newton’s method is either unavailable or computationally intensive. There are various approaches to approximate the Hessian, but the common feature among quasi-Newton methods is that the Hessian approximation must satisfy the secant condition. To be more precise, consider  $\mathbf{s}^t$  and  $\mathbf{y}^t$  as the variable and gradient variations, explicitly defined as

$$(3) \quad \mathbf{s}^t := \mathbf{x}^{t+1} - \mathbf{x}^t, \quad \mathbf{y}^t := \nabla f(\mathbf{x}^{t+1}) - \nabla f(\mathbf{x}^t).$$

Then, given the variable variation  $\mathbf{s}^t$  and gradient variation  $\mathbf{y}^t$ , the Hessian approximation matrix in all quasi-Newton methods must satisfy the secant condition

$$(4) \quad \mathbf{B}^{t+1}\mathbf{s}^t = \mathbf{y}^t.$$

This condition is fundamental in quasi-Newton methods because the exact Hessian  $\nabla^2 f(\mathbf{x}^t)$  satisfies this equality when the iterates  $\mathbf{x}^{t+1}$  and  $\mathbf{x}^t$  are close to each other. If we consider the matrix  $\mathbf{B}^{t+1}$  as the unknown matrix, the system of equations in (4) does not have a unique solution. Different quasi-Newton methods enforce different conditions on the matrix  $\mathbf{B}^{t+1}$  to come up with a unique update. This extra condition is typically a proximity condition that ensures that  $\mathbf{B}^{t+1}$  is close to the previous Hessian approximation matrix  $\mathbf{B}^t$  [4, 22, 12]. In particular, the Broyden-Fletcher-Goldfarb-Shanno (BFGS) method defines the update of Hessian approximation matrix as

$$(5) \quad \mathbf{B}^{t+1} = \mathbf{B}^t + \frac{\mathbf{y}^t \mathbf{y}^{tT}}{\mathbf{y}^{tT} \mathbf{s}^t} - \frac{\mathbf{B}^t \mathbf{s}^t \mathbf{s}^{tT} \mathbf{B}^t}{\mathbf{s}^{tT} \mathbf{B}^t \mathbf{s}^t}.$$

The BFGS method is popular not only for its strong numerical performance relative to the gradient descent method, but also because it is shown to exhibit a superlinear convergence rate [4], thereby providing a theoretical guarantee of superior performance. In fact, it can be shown that, the BFGS update satisfies the condition

$$(6) \quad \lim_{t \rightarrow \infty} \frac{\|(\mathbf{B}^t - \nabla^2 f(\mathbf{x}^*))\mathbf{s}^t\|}{\|\mathbf{s}^t\|} = 0,$$

known as the Dennis-Moré condition, which is both necessary and sufficient for superlinear convergence [12]. This result solidifies quasi-Newton methods as a strong alternative to first order methods when exact second-order information is unavailable. However, implementation of the BFGS method is not feasible when the number of functions  $n$  is large, due to its high computational complexity on the order  $\mathcal{O}(np + p^2)$ . In the following section, we propose a novel incremental BFGS method that has the computational complexity of  $\mathcal{O}(p^2)$  per iteration and converges at a superlinear rate.

**3. IQN: Incremental aggregated BFGS.** We propose an incremental aggregated BFGS algorithm, which we call the Incremental Quasi-Newton (IQN) method. The IQN method is incremental in that, at each iteration, only the information associated with a single function  $f_i$  is updated. The particular function is chosen by cyclicly iterating through the  $n$  functions. The IQN method is aggregated in that the aggregate of the most recently observed information of all functions  $f_1, \dots, f_n$  is used to compute the updated variable  $\mathbf{x}^{t+1}$ .

In the proposed method, we consider  $\mathbf{z}_1^t, \dots, \mathbf{z}_n^t$  as the copies of the variable  $\mathbf{x}$  at time  $t$  associated with the functions  $f_1, \dots, f_n$ , respectively. Likewise, define  $\nabla f_i(\mathbf{z}_i^t)$  as the gradient corresponding to the  $i$ -th function. Further, consider  $\mathbf{B}_i^t$  as a positive definite matrix which approximates the  $i$ -th component Hessian  $\nabla^2 f_i(\mathbf{x}^t)$ . We refer to  $\mathbf{z}_i^t$ ,  $\nabla f_i(\mathbf{z}_i^t)$ , and  $\mathbf{B}_i^t$  as the information corresponding to the  $i$ -th function  $f_i$  at step  $t$ . Note that the functions' information is stored in a shared memory as shown in Fig. 1. To introduce the IQN method, we first explain the mechanism for computing the updated variable  $\mathbf{x}^{t+1}$  using the stored information  $\{\mathbf{z}_i^t, \nabla f_i(\mathbf{z}_i^t), \mathbf{B}_i^t\}_{i=1}^n$ . Then, we elaborate on the scheme for updating the information of the functions.

To derive the full variable update, consider the second order approximation of

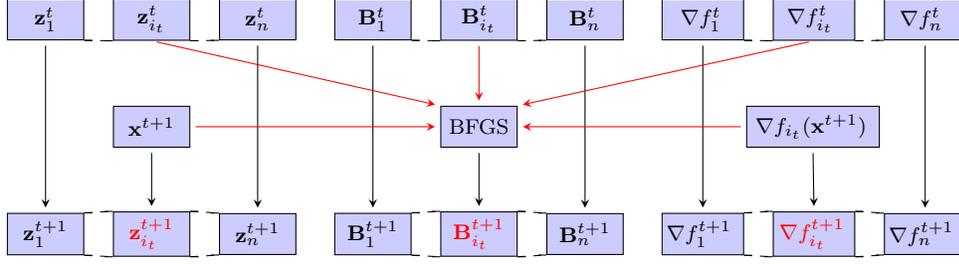


FIG. 1. The updating scheme for variables, gradients, and Hessian approximation matrices of function  $f_{i_t}$  at step  $t$ . The red arrows indicate the terms used in the update of  $\mathbf{B}_{i_t}^{t+1}$  using the BFGS update in (15). The black arrows show the updates of all variables and gradients. The terms  $\mathbf{z}_{i_t}^{t+1}$  and  $\nabla f_{i_t}^{t+1}$  are updated as  $\mathbf{x}^{t+1}$  and  $\nabla f_{i_t}(\mathbf{x}^{t+1})$ , respectively. All others  $\mathbf{z}_j^{t+1}$  and  $\nabla f_j^{t+1}$  are set as  $\mathbf{z}_j^t$  and  $\nabla f_j^t$ , respectively.

the objective function  $f_i(\mathbf{x})$  centered around its current iterate  $\mathbf{z}_i^t$ ,

$$(7) \quad f_i(\mathbf{x}) \approx f_i(\mathbf{z}_i^t) + \nabla f_i(\mathbf{z}_i^t)^T (\mathbf{x} - \mathbf{z}_i^t) + \frac{1}{2} (\mathbf{x} - \mathbf{z}_i^t)^T \nabla^2 f_i(\mathbf{z}_i^t) (\mathbf{x} - \mathbf{z}_i^t).$$

As in traditional quasi-Newton methods, we replace the  $i$ -th Hessian  $\nabla^2 f_i(\mathbf{z}_i^t)$  by  $\mathbf{B}_i^t$ . Using the approximation matrices in place of Hessians, the complete (aggregate) function  $f(\mathbf{x})$  can be approximated with

$$(8) \quad f(\mathbf{x}) \approx \frac{1}{n} \sum_{i=1}^n \left[ f_i(\mathbf{z}_i^t) + \nabla f_i(\mathbf{z}_i^t)^T (\mathbf{x} - \mathbf{z}_i^t) + \frac{1}{2} (\mathbf{x} - \mathbf{z}_i^t)^T \mathbf{B}_i^t (\mathbf{x} - \mathbf{z}_i^t) \right].$$

Note that the right hand side of (8) is a quadratic approximation of the function  $f$  based on the available information at step  $t$ . Hence, the updated iterate  $\mathbf{x}^{t+1}$  can be defined as the minimizer of the quadratic program in (8), explicitly given by

$$(9) \quad \mathbf{x}^{t+1} = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t \right)^{-1} \left[ \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t \mathbf{z}_i^t - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{z}_i^t) \right].$$

First note that the update in (9) shows that the updated variable  $\mathbf{x}^{t+1}$  is a function of the stored information of all functions  $f_1, \dots, f_n$ . Furthermore, we use the aggregated information of variables, gradients, and the quasi-Newton Hessian approximations to evaluate the updated variable. This is done to vanish the noise in approximating both gradients and Hessians as the sequence approaches the optimal argument.

*Remark 1.* Given the BFGS Hessian approximation matrices  $\{\mathbf{B}_i^t\}_{i=1}^n$  and gradients  $\{\nabla f_i(\mathbf{z}_i^t)\}_{i=1}^n$ , one may consider an update more akin to traditional descent-based methods, i.e.

$$(10) \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \left( \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t \right)^{-1} \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{z}_i^t).$$

To evaluate the advantage of the proposed update for IQN in (9) relative to the update in (10), we proceed to study the Taylor's expansion that leads to the update in (10).

It can be shown that the update in (10) is the outcome of the following approximation

$$(11) \quad f(\mathbf{x}) \approx \frac{1}{n} \sum_{i=1}^n \left[ f_i(\mathbf{z}_i^t) + \nabla f_i(\mathbf{z}_i^t)^T (\mathbf{x} - \mathbf{z}_i^t) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^t)^T \mathbf{B}_i^t (\mathbf{x} - \mathbf{x}^t) \right].$$

Observe that the linear term in (11) is centered at  $\mathbf{z}_i^t$ , while the quadratic term is approximated near the iterate  $\mathbf{x}^t$ . This inconsistency in the Taylor's expansion of each function  $f_i$  leads to an inaccurate second-order approximation, and subsequently a slower incremental quasi-Newton method.

Thus far we have discussed the procedure to compute the updated variable  $\mathbf{x}^{t+1}$  given the local iterates, gradients, and Hessian approximations at time  $t$ . Now, it remains to show how we update the local information of functions  $f_1, \dots, f_n$  using the variable  $\mathbf{x}^{t+1}$ . In each iteration of the IQN method, we update the local information of only a single function, chosen in a cyclic manner. Defining  $i_t$  to be the index of the function selected at time  $t$ , we update the local variables  $\mathbf{z}_{i_t}^{t+1}$ ,  $\nabla f_{i_t}(\mathbf{z}_{i_t}^{t+1})$ , and  $\mathbf{B}_{i_t}^{t+1}$  using the updated variable  $\mathbf{x}^{t+1}$  while all other local variables remain unchanged. In particular, the variables  $\mathbf{z}_i$  are updated as

$$(12) \quad \mathbf{z}_{i_t}^{t+1} = \mathbf{x}^{t+1}, \quad \mathbf{z}_i^{t+1} = \mathbf{z}_i^t \quad \text{for all } i \neq i_t.$$

Observe in the update in (12) that the variable associated with the function  $f_{i_t}$  is set to be the updated variable  $\mathbf{x}^{t+1}$  while the other iterates are simply kept as their previous value. Likewise, we update the table of gradients accordingly with the gradient of  $f_{i_t}$  evaluated at the new variable  $\mathbf{x}^{t+1}$ . The rest of gradients stored in the memory will stay unchanged, i.e.,

$$(13) \quad \nabla f_{i_t}(\mathbf{z}_{i_t}^{t+1}) = \nabla f_{i_t}(\mathbf{x}^{t+1}), \quad \nabla f_i(\mathbf{z}_i^{t+1}) = \nabla f_i(\mathbf{z}_i^t) \quad \text{for all } i \neq i_t.$$

To update the curvature information, it would be ideal to compute the Hessian  $\nabla^2 f_{i_t}(\mathbf{x}^{t+1})$  and update the curvature information following the schemes for variables in (12) and gradients in (13). However, our focus is on the applications that the computation of the Hessian is either impossible or computationally expensive. Hence, to the update curvature approximation matrix  $\mathbf{B}_{i_t}^t$  corresponding to the function  $f_{i_t}$ , we use the steps of BFGS in (5). To do so, we define variable and gradient variations associated with each individual function  $f_i$  as

$$(14) \quad \mathbf{s}_i^t := \mathbf{z}_i^{t+1} - \mathbf{z}_i^t, \quad \mathbf{y}_i^t := \nabla f_i(\mathbf{z}_i^{t+1}) - \nabla f_i(\mathbf{z}_i^t),$$

respectively. The Hessian approximation  $\mathbf{B}_{i_t}^t$  corresponding to the function  $f_{i_t}$  can be computed using the update of BFGS as

$$(15) \quad \mathbf{B}_i^{t+1} = \mathbf{B}_i^t + \frac{\mathbf{y}_i^t \mathbf{y}_i^{tT}}{\mathbf{y}_i^{tT} \mathbf{s}_i^t} - \frac{\mathbf{B}_i^t \mathbf{s}_i^t \mathbf{s}_i^{tT} \mathbf{B}_i^t}{\mathbf{s}_i^{tT} \mathbf{B}_i^t \mathbf{s}_i^t}, \quad \text{for } i = i_t.$$

Again, the Hessian approximation matrices for all other functions remain unchanged, i.e.,  $\mathbf{B}_i^{t+1} = \mathbf{B}_i^t$  for  $i \neq i_t$ . The system of updates in (12)-(15) explains the mechanism of updating the information of the function  $f_{i_t}$  at step  $t$ . Notice that to update the Hessian approximation matrix for the  $i_t$ -th function there is no need to store the variations in (14), since the old variables  $\mathbf{z}_i^t$  and  $\nabla f_i(\mathbf{z}_i^t)$  are available in memory and the updated versions  $\mathbf{z}_i^{t+1} = \mathbf{x}^{t+1}$  and  $\nabla f_i(\mathbf{z}_i^{t+1}) = \nabla f_i(\mathbf{x}^{t+1})$  are evaluated at step  $t$ ; see Fig. 1 for more details.

Because of the cyclic update scheme, the set of iterates  $\{\mathbf{z}_1^t, \mathbf{z}_2^t, \dots, \mathbf{z}_n^t\}$  is equal to the set  $\{\mathbf{x}^t, \mathbf{x}^{t-1}, \dots, \mathbf{x}^{t-n+1}\}$ , and, therefore, the set of variables used in the update of IQN is the set of the last  $n$  iterates. The update of IQN in (9) incorporates the information of all the functions  $f_1, \dots, f_n$  to compute the updated variable  $\mathbf{x}^{t+1}$ ; however, it uses delayed variables, gradients, and Hessian approximations rather than the the updated variable  $\mathbf{x}^{t+1}$  for all functions as in classic quasi-Newton methods. The use of delay allows IQN to update the information of a single function at each iteration, thus reducing the computational complexity relative to classic quasi-Newton methods.

Although the update in (9) is helpful in understanding the rationale behind the IQN method, it cannot be implemented at a low computation cost, since it requires computation of the sums  $\sum_{i=1}^n \mathbf{B}_i^t$ ,  $\sum_{i=1}^n \mathbf{B}_i^t \mathbf{z}_i^t$ , and  $\sum_{i=1}^n \nabla f_i(\mathbf{z}_i^t)$  as well as computing the inversion  $(\sum_{i=1}^n \mathbf{B}_i^t)^{-1}$ . In the following section, we introduce an efficient implementation of the IQN method that has the computational complexity of  $\mathcal{O}(p^2)$ .

**3.1. Efficient implementation of IQN.** To see that the updating scheme in (9) requires evaluation of only a single gradient and Hessian approximation matrix per iteration, consider writing the update as

$$(16) \quad \mathbf{x}^{t+1} = (\tilde{\mathbf{B}}^t)^{-1} (\mathbf{u}^t - \mathbf{g}^t),$$

where we define  $\tilde{\mathbf{B}}^t := \sum_{i=1}^n \mathbf{B}_i^t$  as the aggregate Hessian approximation,  $\mathbf{u}^t := \sum_{i=1}^n \mathbf{B}_i^t \mathbf{z}_i^t$  as the aggregate Hessian-variable product, and  $\mathbf{g}^t := \sum_{i=1}^n \nabla f_i(\mathbf{z}_i^t)$  as the aggregate gradient. Then, given that at step  $t$  only a single index  $i_t$  is updated, we can evaluate these variables for step  $t+1$  as

$$(17) \quad \tilde{\mathbf{B}}^{t+1} = \tilde{\mathbf{B}}^t + (\mathbf{B}_{i_t}^{t+1} - \mathbf{B}_{i_t}^t),$$

$$(18) \quad \mathbf{u}^{t+1} = \mathbf{u}^t + (\mathbf{B}_{i_t}^{t+1} \mathbf{z}_{i_t}^{t+1} - \mathbf{B}_{i_t}^t \mathbf{z}_{i_t}^t),$$

$$(19) \quad \mathbf{g}^{t+1} = \mathbf{g}^t + (\nabla f_{i_t}(\mathbf{z}_{i_t}^{t+1}) - \nabla f_{i_t}(\mathbf{z}_{i_t}^t)).$$

Thus, only  $\mathbf{B}_{i_t}^{t+1}$  and  $\nabla f_{i_t}(\mathbf{z}_{i_t}^{t+1})$  are required to be computed at step  $t$ .

Although the updates in (17)-(19) have low computational complexity, the update in (16) requires computing  $(\tilde{\mathbf{B}}^t)^{-1}$  which has a computational complexity of  $\mathcal{O}(p^3)$ . This inversion can be avoided by simplifying the update in (17) as

$$(20) \quad \tilde{\mathbf{B}}^{t+1} = \tilde{\mathbf{B}}^t + \frac{\mathbf{y}_{i_t}^t \mathbf{y}_{i_t}^{tT}}{\mathbf{y}_{i_t}^{tT} \mathbf{s}_{i_t}^t} - \frac{\mathbf{B}_{i_t}^t \mathbf{s}_{i_t}^t \mathbf{s}_{i_t}^{tT} \mathbf{B}_{i_t}^t}{\mathbf{s}_{i_t}^{tT} \mathbf{B}_{i_t}^t \mathbf{s}_{i_t}^t}.$$

To derive the expression in (20) we have substituted the difference  $\mathbf{B}_{i_t}^{t+1} - \mathbf{B}_{i_t}^t$  by its rank two expression in (15). Given the matrix  $(\tilde{\mathbf{B}}^t)^{-1}$ , by applying the Sherman-Morrison formula twice to the update in (20) we can compute  $(\tilde{\mathbf{B}}^{t+1})^{-1}$  as

$$(21) \quad (\tilde{\mathbf{B}}^{t+1})^{-1} = \mathbf{U}^t + \frac{\mathbf{U}^t (\mathbf{B}_{i_t}^t \mathbf{s}_{i_t}^t) (\mathbf{B}_{i_t}^t \mathbf{s}_{i_t}^t)^T \mathbf{U}^t}{\mathbf{s}_{i_t}^{tT} \mathbf{B}_{i_t}^t \mathbf{s}_{i_t}^t - (\mathbf{B}_{i_t}^t \mathbf{s}_{i_t}^t)^T \mathbf{U}^t (\mathbf{B}_{i_t}^t \mathbf{s}_{i_t}^t)},$$

where the matrix  $\mathbf{U}^t$  is evaluated as

$$(22) \quad \mathbf{U}^t = (\tilde{\mathbf{B}}^t)^{-1} - \frac{(\tilde{\mathbf{B}}^t)^{-1} \mathbf{y}_{i_t}^t \mathbf{y}_{i_t}^{tT} (\tilde{\mathbf{B}}^t)^{-1}}{\mathbf{y}_{i_t}^{tT} \mathbf{s}_{i_t}^t + \mathbf{y}_{i_t}^{tT} (\tilde{\mathbf{B}}^t)^{-1} \mathbf{y}_{i_t}^t}.$$

**Algorithm 1** Incremental Quasi-Newton (IQN) method

---

**Require:**  $\mathbf{x}^0, \{\nabla f_i(\mathbf{x}^0)\}_{i=1}^n, \{\mathbf{B}_i^0\}_{i=1}^n$   
1: Set  $\mathbf{z}_1^0 = \dots = \mathbf{z}_n^0 = \mathbf{x}^0$   
2: Set  $(\tilde{\mathbf{B}}^0)^{-1} = (\sum_{i=1}^n \mathbf{B}_i^0)^{-1}$ ,  $\mathbf{u}^0 = \sum_{i=1}^n \mathbf{B}_i^0 \mathbf{x}^0$ ,  $\mathbf{g}^0 = \sum_{i=1}^n \nabla f_i(\mathbf{x}^0)$   
3: **for**  $t = 0, 1, 2, \dots$  **do**  
4:   Set  $i_t = (t \bmod n) + 1$   
5:   Compute  $\mathbf{x}^{t+1} = (\tilde{\mathbf{B}}^t)^{-1} (\mathbf{u}^t - \mathbf{g}^t)$  [cf. (16)]  
6:   Compute  $\mathbf{s}_{i_t}^{t+1}$ ,  $\mathbf{y}_{i_t}^{t+1}$  [cf. (14)], and  $\mathbf{B}_{i_t}^{t+1}$  [cf. (15)]  
7:   Update  $\mathbf{u}^{t+1}$  [cf. (18)],  $\mathbf{g}^{t+1}$  [cf. (19)], and  $(\tilde{\mathbf{B}}^{t+1})^{-1}$  [cf. (21), (22)]  
8:   Update the functions' information tables as in (12), (13), and (15)  
9: **end for**

---

The computational complexity of the updates in (21) and (22) is of the order  $\mathcal{O}(p^2)$  rather than the  $\mathcal{O}(p^3)$  cost of computing the inverse directly. Therefore, the overall cost of IQN is of the order  $\mathcal{O}(p^2)$  which is substantially lower than  $\mathcal{O}(np^2)$  of deterministic quasi-Newton methods.

The complete IQN algorithm is outlined in Algorithm 1. Beginning with initial variable  $\mathbf{x}^0$  and gradient and Hessian estimates  $\nabla f_i(\mathbf{x}^0)$  and  $\mathbf{B}_i^0$  for all  $i$ , each variable copy  $\mathbf{z}_i^0$  is set to  $\mathbf{x}^0$  in Step 1 and initial values are set for  $\mathbf{u}^0$ ,  $\mathbf{g}^0$  and  $(\tilde{\mathbf{B}}^0)^{-1}$  in Step 2. For all  $t$ , in Step 4 the index  $i_t$  of the next function to update is selected cyclically. The variable  $\mathbf{x}^{t+1}$  is computed according to the update in (16) in Step 5. In Step 6, the variable  $\mathbf{s}_{i_t}^{t+1}$  and gradient  $\mathbf{y}_{i_t}^{t+1}$  variations are evaluated as in (14) to compute the BFGS matrix  $\mathbf{B}_{i_t}^{t+1}$  from the update in (15). This information, as well as the updated variable and its gradient, are used in Step 7 to update  $\mathbf{u}^{t+1}$  and  $\mathbf{g}^{t+1}$  as in (18) and (19), respectively. The inverse matrix  $(\tilde{\mathbf{B}}^{t+1})^{-1}$  is also computed by following the expressions in (21) and (22). Finally in Step 8, we update the variable, gradient, and Hessian approximation tables based on the policies in (12), (13), and (15), respectively.

**4. Convergence Analysis.** In this section, we study the convergence rate of the proposed IQN method. We first establish its local linear convergence rate, then demonstrate limit properties of the Hessian approximations, and finally show that in a region local to the optimal point the sequence of residuals converges at a superlinear rate. To prove these results we make two main assumptions, both of which are standard in the analysis of quasi-Newton methods.

**Assumption 1.** *There exist positive constants  $0 < \mu \leq L$  such that, for all  $i$  and  $\mathbf{x}, \hat{\mathbf{x}} \in \mathbb{R}^p$ , we can write*

$$(23) \quad \mu \|\mathbf{x} - \hat{\mathbf{x}}\|^2 \leq (\nabla f_i(\mathbf{x}) - \nabla f_i(\hat{\mathbf{x}}))^T (\mathbf{x} - \hat{\mathbf{x}}) \leq L \|\mathbf{x} - \hat{\mathbf{x}}\|^2.$$

**Assumption 2.** *There exists a positive constant  $0 < \tilde{L}$  such that, for all  $i$  and  $\mathbf{x}, \hat{\mathbf{x}} \in \mathbb{R}^p$ , we can write*

$$(24) \quad \|\nabla^2 f_i(\mathbf{x}) - \nabla^2 f_i(\hat{\mathbf{x}})\| \leq \tilde{L} \|\mathbf{x} - \hat{\mathbf{x}}\|.$$

The lower bound in (23) implies that the functions  $f_i$  are strongly convex with constant  $\mu$ , and the upper bound shows that the gradients  $\nabla f_i$  are Lipschitz continuous with parameter  $L$ .

The condition in Assumption 2, states that the Hessians  $\nabla^2 f_i$  are Lipschitz continuous with constant  $\tilde{L}$ . This assumption is commonly made in the analyses of Newton's

method [21] and quasi-Newton algorithms [4, 22, 12]. According to Lemma 3.1 in [4], Lipschitz continuity of the Hessians with constant  $\tilde{L}$  implies that for  $i = 1, \dots, n$  and arbitrary vectors  $\mathbf{x}, \tilde{\mathbf{x}}, \hat{\mathbf{x}} \in \mathbb{R}^p$  we can write

$$(25) \quad \|\nabla^2 f_i(\tilde{\mathbf{x}})(\mathbf{x} - \hat{\mathbf{x}}) - (\nabla f_i(\mathbf{x}) - \nabla f_i(\hat{\mathbf{x}}))\| \leq \tilde{L} \|\mathbf{x} - \hat{\mathbf{x}}\| \max\{\|\mathbf{x} - \tilde{\mathbf{x}}\|, \|\hat{\mathbf{x}} - \tilde{\mathbf{x}}\|\}.$$

We use the inequality in (25) in the process of proving the convergence of IQN.

The goal of BFGS quasi-Newton methods is to approximate the objective function Hessian using the first-order information. Likewise, in the incremental BFGS method, we aim to show that the Hessian approximation matrices for all functions  $f_1, \dots, f_n$  are close to the exact Hessian. In the following lemma, we study the difference between the  $i$ -th optimal Hessian  $\nabla^2 f_i(\mathbf{x}^*)$  and its approximation  $\mathbf{B}_i^t$  over time.

LEMMA 1. *Consider the proposed IQN method in (9). Further, let  $i$  be the index of the updated function at step  $t$ , i.e.,  $i = i_t$ . Define the residual sequence for function  $f_i$  as  $\sigma_i^t := \max\{\|\mathbf{z}_i^{t+1} - \mathbf{x}^*\|, \|\mathbf{z}_i^t - \mathbf{x}^*\|\}$  and set  $\mathbf{M} = \nabla^2 f_i(\mathbf{x}^*)^{-1/2}$ . If Assumptions 1 and 2 hold and the condition  $\sigma_i^t < m/(3\tilde{L})$  is satisfied then*

$$(26) \quad \|\mathbf{B}_i^{t+1} - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} \leq \left[ (1 - \alpha\theta_i^{t^2})^{1/2} + \alpha_3\sigma_i^t \right] \|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} + \alpha_4\sigma_i^t,$$

where  $\alpha, \alpha_3$ , and  $\alpha_4$  are some positive bounded constants and

$$(27) \quad \theta_i^t = \frac{\|\mathbf{M}(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*))\mathbf{s}_i^t\|}{\|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} \|\mathbf{M}^{-1}\mathbf{s}_i^t\|} \quad \text{for } \mathbf{B}_i^t \neq \nabla^2 f_i(\mathbf{x}^*), \quad \theta_i^t = 0 \quad \text{for } \mathbf{B}_i^t = \nabla^2 f_i(\mathbf{x}^*).$$

*Proof.* See Appendix A. □

The result in (26) establishes an upper bound for the weighted norm  $\|\mathbf{B}_i^{t+1} - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}}$  with respect to its previous value  $\|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}}$  and the sequence  $\sigma_i^t := \max\{\|\mathbf{z}_i^{t+1} - \mathbf{x}^*\|, \|\mathbf{z}_i^t - \mathbf{x}^*\|\}$ , when the variables are in a neighborhood of the optimal solution such that  $\sigma_i^t < m/(3\tilde{L})$ . Indeed, the result in (26) holds only for the index  $i = i_t$  and for the rest of indices we have  $\|\mathbf{B}_i^{t+1} - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} = \|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}}$  simply by definition of the cyclic update. Note that if the residual sequence  $\sigma_i^t$  associated with  $f_i$  approaches zero, we can simplify (26) as

$$(28) \quad \|\mathbf{B}_i^{t+1} - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} \lesssim (1 - \alpha\theta_i^{t^2})^{1/2} \|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}}.$$

The equation in (28) implies that if  $\theta_i^t$  is always strictly larger than zero, the sequence  $\|\mathbf{B}_i^{t+1} - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}}$  approaches zero. If not, then the sequence  $\theta_i^t$  converges to zero which implies the Dennis-Moré condition from (6), i.e.

$$(29) \quad \lim_{t \rightarrow \infty} \frac{\|(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*))\mathbf{s}_i^t\|}{\|\mathbf{s}_i^t\|} = 0.$$

Therefore, under both conditions the result in (29) holds. This is true since the limit  $\lim_{t \rightarrow \infty} \|\mathbf{B}_i^{t+1} - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} = 0$  yields the result in (29).

Based on this intuition, we proceed to show that the sequence  $\sigma_i^t$  converges to zero for all  $i = 1, \dots, n$ . To do so, we show that the sequence  $\|\mathbf{z}_i^t - \mathbf{x}^*\|$  is linearly convergent for all  $i = 1, \dots, n$ . To achieve this goal we first prove an upper bound for the error  $\|\mathbf{x}^{t+1} - \mathbf{x}^*\|$  of IQN in the following lemma.

LEMMA 2. Consider the proposed IQN method in (9). If the conditions in Assumptions 1 and 2 hold, then the sequence of iterates generated by IQN satisfies

$$(30) \quad \|\mathbf{x}^{t+1} - \mathbf{x}^*\| \leq \frac{\tilde{L}\Gamma^t}{n} \sum_{i=1}^n \|\mathbf{z}_i^t - \mathbf{x}^*\|^2 + \frac{\Gamma^t}{n} \sum_{i=1}^n \|(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)) (\mathbf{z}_i^t - \mathbf{x}^*)\|,$$

where  $\Gamma^t := \|((1/n) \sum_{i=1}^n \mathbf{B}_i^t)^{-1}\|$ .

*Proof.* See Appendix B.  $\square$

Lemma 2 shows that the residual  $\|\mathbf{x}^{t+1} - \mathbf{x}^*\|$  is bounded above by a sum of quadratic and linear terms of the last  $n$  residuals. This can eventually lead to a superlinear convergence rate by establishing the linear term converges to zero at a fast rate, leaving us with an upper bound of quadratic terms only. First, however, we establish a local linear convergence rate in the proceeding theorem to show that the sequence  $\sigma_i^t$  converges to zero.

LEMMA 3. Consider the proposed IQN method in (9). If Assumptions 1 and 2 hold, then, for any  $r \in (0, 1)$  there are positive constants  $\epsilon(r)$  and  $\delta(r)$  such that if  $\|\mathbf{x}^0 - \mathbf{x}^*\| < \epsilon(r)$  and  $\|\mathbf{B}_i^0 - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} < \delta(r)$  for  $\mathbf{M} = \nabla^2 f_i(\mathbf{x}^*)^{-1/2}$  and  $i = 1, 2, \dots, n$ , the sequence of iterates generated by IQN satisfies

$$(31) \quad \|\mathbf{x}^t - \mathbf{x}^*\| \leq r^{\lfloor \frac{t-1}{n} \rfloor + 1} \|\mathbf{x}^0 - \mathbf{x}^*\|.$$

Moreover, the sequences of norms  $\{\|\mathbf{B}_i^t\|\}$  and  $\{\|(\mathbf{B}_i^t)^{-1}\|\}$  are uniformly bounded.

*Proof.* See Appendix C.  $\square$

The result in Lemma 3 shows that the sequence of iterates generated by IQN has a local linear convergence rate after each pass over all functions. Consequently, we obtain that the  $i$ -th residual sequence  $\sigma_i^t$  is linearly convergent for all  $i$ . Note that Lemma 3 can be considered as an extension of Theorem 3.2 in [4] for incremental settings. Following the arguments in (28) and (29), we use the summability of the sequence  $\sigma_i^t$  along with the result in Lemma 1 to prove Dennis-Moré condition for all functions  $f_i$ .

PROPOSITION 4. Consider the proposed IQN method in (9). Assume that the hypotheses in Lemmata 1 and 3 are satisfied. Then, for all  $i = 1, \dots, n$  it holds,

$$(32) \quad \lim_{t \rightarrow \infty} \frac{\|(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)) \mathbf{s}_i^t\|}{\|\mathbf{s}_i^t\|} = 0.$$

*Proof.* See Appendix D.  $\square$

The statement in Proposition 4 indicates that for each function  $f_i$  the Dennis-Moré condition holds. In the tradition quasi-Newton methods the Dennis-Moré condition is sufficient to show that the method is superlinearly convergent. However, the same argument does not hold for the proposed IQN method, since we can't recover the Dennis-Moré condition for the global objective function  $f$  from the result in Proposition 4. In other words, the result in (32) does not imply the limit in (6) required in the superlinear convergence analysis of quasi-Newton methods. Therefore, here we pursue a different approach and seek to prove that the linear terms

$(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*))(\mathbf{z}_i^t - \mathbf{x}^*)$  in (30) converge to zero at a superlinear rate, i.e., for all  $i$  we can write  $\lim_{t \rightarrow \infty} \|(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*))(\mathbf{z}_i^t - \mathbf{x}^*)\| / \|\mathbf{z}_i^t - \mathbf{x}^*\| = 0$ . If we establish this result, it follows from the result in Lemma 2 that the sequence of residuals  $\|\mathbf{x}^t - \mathbf{x}^*\|$  converges to zero superlinearly.

We continue the analysis of the proposed IQN method by establishing a generalized limit property that follows from the Dennis-Moré criterion in (6). In the following lemma, we leverage the local linear convergence of the iterates  $\mathbf{x}^t$  to show that the vector  $\mathbf{z}_i^t - \mathbf{x}^*$  lies in the null space of  $\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)$  as  $t$  approaches infinity.

LEMMA 5. *Consider the proposed IQN method in (9). Assume that the hypotheses in Lemmata 1 and 3 are satisfied. As  $t$  approaches infinity, the following holds for all  $i$ ,*

$$(33) \quad \lim_{t \rightarrow \infty} \frac{\|(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*))(\mathbf{z}_i^t - \mathbf{x}^*)\|}{\|\mathbf{z}_i^t - \mathbf{x}^*\|} = 0.$$

*Proof.* See Appendix E. □

The result in Lemma 5 can thus be used in conjunction with Lemma 2 to show that the residual  $\|\mathbf{x}^{t+1} - \mathbf{x}^*\|$  is bounded by a sum of quadratic terms of previous residuals and a term that converges to zero superlinearly. This result leads us to the following result, namely the local superlinear convergence of the sequence of residuals with respect to the average sequence, stated in the following theorem.

THEOREM 6. *Consider the proposed IQN method in (9). Suppose that the conditions in the hypotheses of Lemmata 1 and 3 are valid. Then, the sequence of residuals  $\|\mathbf{x}^t - \mathbf{x}^*\|$  satisfies*

$$(34) \quad \lim_{t \rightarrow \infty} \frac{\|\mathbf{x}^t - \mathbf{x}^*\|}{\frac{1}{n}(\|\mathbf{x}^{t-1} - \mathbf{x}^*\| + \dots + \|\mathbf{x}^{t-n} - \mathbf{x}^*\|)} = 0.$$

*Proof.* See Appendix F. □

The result in (34) shows a mean-superlinear convergence rate for the sequence of iterates generated by IQN. To be more precise, it shows that the ratio that captures the error at step  $t$  divided by the average of last  $n$  errors converges to zero. This is not equivalent to the classic Q-superlinear convergence for full-batch quasi-Newton methods, i.e.,  $\lim_{t \rightarrow \infty} \|\mathbf{x}^{t+1} - \mathbf{x}^*\| / \|\mathbf{x}^t - \mathbf{x}^*\| = 0$ . Although Q-superlinear convergence of the residuals  $\|\mathbf{x}^t - \mathbf{x}^*\|$  is not provable, we can show that there exists a subsequence of the sequence  $\|\mathbf{x}^t - \mathbf{x}^*\|$  that converges to zero superlinearly. In addition, there exists a superlinearly convergent sequence that is an upper bound for the original sequence of errors  $\|\mathbf{x}^t - \mathbf{x}^*\|$ . We formalize these results in the following theorem.

THEOREM 7. *Consider the proposed IQN method in (9). Suppose that the conditions in the hypotheses of Lemmata 1 and 3 are valid. Then, there exists a subsequence of  $\|\mathbf{x}^t - \mathbf{x}^*\|$  that converges to zero superlinearly. Moreover, there exists a sequence  $\zeta^t$  such that  $\|\mathbf{x}^t - \mathbf{x}^*\| \leq \zeta^t$  for all  $t \geq 0$ , and the sequence  $\zeta^t$  converges to zero at a superlinear rate, i.e.,*

$$(35) \quad \lim_{t \rightarrow \infty} \frac{\zeta^{t+1}}{\zeta^t} = 0.$$

*Proof.* See Appendix G.  $\square$

The first result in Theorem 7 states that although the whole sequence  $\|\mathbf{x}^t - \mathbf{x}^*\|$  is not necessarily superlinearly convergent, there exists a subsequence of the sequence  $\|\mathbf{x}^t - \mathbf{x}^*\|$  that converges at a superlinear rate. The second claim in Theorem 7 establishes R-superlinear convergence rate of the whole sequence  $\|\mathbf{x}^t - \mathbf{x}^*\|$ . In other words, it guarantees that  $\|\mathbf{x}^t - \mathbf{x}^*\|$  is upper bounded by a superlinearly convergent sequence.

**5. Related Works.** Various methods have been studied in the literature to improve the performance of traditional full-batch optimization algorithms. The most famous method for reducing the computational complexity of gradient descent (GD) is stochastic gradient descent (SGD), which uses the gradient of a single randomly chosen function to approximate the full-gradient [2]. Incremental gradient descent method (IGD) is similar to SGD except the function is chosen in a cyclic routine [1]. Both SGD and IGD suffer from slow sublinear convergence rate because of the noise of gradient approximation. The incremental aggregated methods, which use memory to aggregate the gradients of all  $n$  functions, are successful in reducing the noise of gradient approximation to achieve linear convergence rate [26, 28, 9, 13]. The work in [26] suggests a random selection of functions which leads to stochastic average gradient method (SAG), while the works in [1, 11, 17] use a cyclic scheme.

Moving beyond first order information, there have been stochastic quasi-Newton methods to approximate Hessian information [29, 18, 19, 20, 10]. All of these stochastic quasi-Newton methods reduce computational cost of quasi-Newton methods by updating only a randomly chosen single or small subset of gradients at each iteration. However, they are not able to recover the superlinear convergence rate of quasi-Newton methods [4, 22, 12]. The incremental Newton method (NIM) in [25] is the only incremental method shown to have a superlinear convergence rate; however, the Hessian function is not always available or computationally feasible. Moreover, the implementation of NIM requires computation of the incremental aggregated Hessian inverse which has the computational complexity of the order  $\mathcal{O}(p^3)$ .

**6. Numerical Results.** We proceed by simulating the performance of IQN on a variety of machine learning problems on both artificial and real datasets. We compare the performance of IQN against a collection of well known first order stochastic and incremental algorithms—namely SAG, SAGA, and IAG. To begin, we look at a simple quadratic program, also equivalent to the solution of linear least squares estimation problem. Consider the objective function to be minimized,

$$(36) \quad \mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^p}{\operatorname{argmin}} f(\mathbf{x}) := \underset{\mathbf{x} \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \mathbf{x}^T \mathbf{A}_i \mathbf{x} + \mathbf{b}_i^T \mathbf{x}.$$

We generate  $\mathbf{A}_i \in \mathbb{R}^{p \times p}$  as a random positive definite matrix and  $\mathbf{b}_i \in \mathbb{R}^p$  as a random vector for all  $i$ . In particular we set the matrices  $\mathbf{A}_i := \operatorname{diag}\{\mathbf{a}_i\}$  and generate random vectors  $\mathbf{a}_i$  with the first  $p/2$  elements chosen from  $[1, 10^{\xi/2}]$  and last  $p/2$  elements chosen from  $[10^{-\xi/2}, 1]$ . The parameter  $\xi$  is used to manually set the condition number for the quadratic program in (36), ranging from  $\xi = 1$  (i.e. small condition number  $10^2$ ) and  $\xi = 2$  (i.e. large condition number  $10^4$ ). The vectors  $\mathbf{b}_i$  are chosen uniformly and randomly from the box  $[0, 10^3]^p$ . The variable dimension is set to be  $p = 10$  and number of functions  $n = 1000$ . Given that we focus on local convergence, we use a constant step size of  $\eta = 1$  for the proposed IQN method while choosing the largest step size allowable by the other methods to converge.

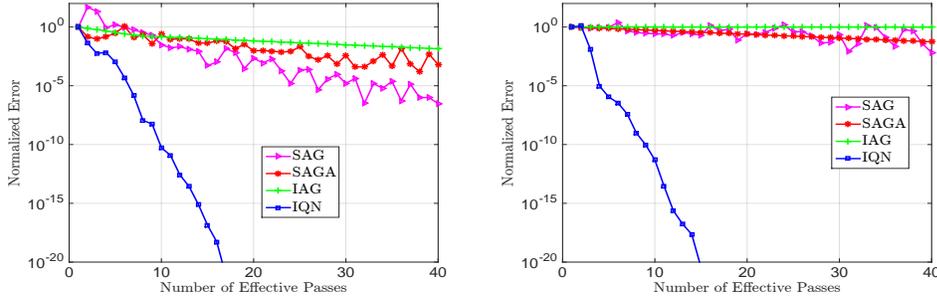


FIG. 2. Convergence results of proposed IQN method in comparison to SAG, SAGA, and IAG. In the left image, we present a sample convergence path of the normalized error on the quadratic program with a small condition number. In the right image, we show the convergence path for the quadratic program with a large condition number. In all cases, IQN provides significant improvement over first order methods, with the difference increasing for larger condition number.

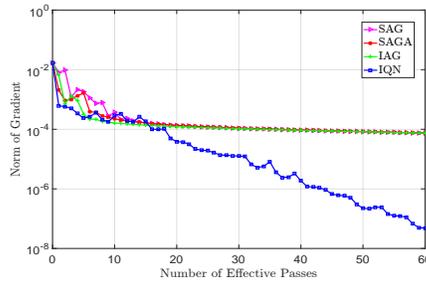


FIG. 3. Convergence results for a sample convergence path for the logistic regression problem on classifying handwritten digits. IQN substantially outperforms the first order methods.

In Figure 2 we present a simulation of the convergence path of the normalized error  $\|\mathbf{x}^t - \mathbf{x}^*\|/\|\mathbf{x}^0 - \mathbf{x}^*\|$  for the quadratic program. In the the left image, we show a sample simulation path for all methods on the quadratic problem with a small condition number. Step sizes of  $\eta = 5 \times 10^{-5}$ ,  $\eta = 10^{-4}$  and  $\eta = 10^{-6}$  were used for SAG, SAGA, and IAG, respectively. These step sizes are tuned to compare the best performance of these methods with IQN. The proposed method reaches a error of  $10^{-10}$  after 10 passes through the data. Alternatively, SAGA achieves the same error of  $10^{-5}$  after 30 passes, while SAG and IAG do not reach  $10^{-5}$  after 40 passes.

In the right image of Figure 2, we repeat the same simulation but with larger condition number. In this case, SAG uses stepsize  $\eta = 2 \times 10^{-4}$  while others remain the same. Observe that while the performance of IQN does not degrade with larger condition number, the first order methods all suffer large degradation. SAG, SAGA, and IAG reach after 40 passes a normalized error of  $6.5 \times 10^{-3}$ ,  $5.5 \times 10^{-2}$ , and  $9.6 \times 10^{-1}$ , respectively. It can be seen that IQN significantly outperforms the first order method for both condition number sizes, with the outperformance increasing for larger condition number. This is an expected result, as first order methods often do not perform well for ill conditioned problems.

**6.1. Logistic regression.** We proceed to numerically evaluate the performance of IQN relative to existing methods on the classification of handwritten digits in the MNIST database [14]. In particular, we solve the binary logistic regression problem.

A logistic regression takes as inputs  $n$  training feature vectors  $\mathbf{u}_i \in \mathbb{R}^p$  with associated labels  $v_i \in \{-1, 1\}$  and outputs a linear classifier  $\mathbf{x}$  to predict the label of unknown feature vectors. For the digit classification problem, each feature vector  $\mathbf{u}_i$  represents a vectorized image and label  $v_i$  its label as one of two digits. We evaluate for any training sample  $i$  the probability of a label  $v_i = 1$  given image  $\mathbf{u}_i$  as  $P(v = 1 | \mathbf{u}) = 1/(1 + \exp(-\mathbf{u}^T \mathbf{x}))$ . The classifier  $\mathbf{x}$  is chosen to be the vector which maximizes the log likelihood across all  $n$  samples. Given  $n$  images  $\mathbf{u}_i$  with associated labels  $v_i$ , the optimization problem for logistic regression is written as

$$(37) \quad \mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^p}{\operatorname{argmin}} f(\mathbf{x}) := \underset{\mathbf{x} \in \mathbb{R}^p}{\operatorname{argmin}} \frac{\lambda}{2} \|\mathbf{x}\|^2 + \frac{1}{n} \sum_{i=1}^n \log[1 + \exp(-v_i \mathbf{u}_i^T \mathbf{x})],$$

where the first term is a regularization term parametrized by  $\lambda \geq 0$ .

For our simulations we select from the MNIST dataset  $n = 1000$  images with dimension  $p = 784$  labelled as one of the digits “0” or “8” and fix the regularization parameter as  $\lambda = 1/n$  and stepsize  $\eta = 0.01$  for all first order methods. In Figure 3 we present the convergence path of IQN relative to existing methods in terms of the norm of the gradient. As in the case of the quadratic program, the IQN performs all gradient-based methods. IQN reaches a gradient magnitude of  $4.8 \times 10^{-8}$  after 60 passes through the data while the SAGA reaches only a magnitude of  $7.4 \times 10^{-5}$  (all other methods perform even worse). Further note that while the first order methods begin to level out after 60 passes, the IQN method continues to descend. These results demonstrate the effectiveness of IQN on a practical machine learning problem with real world data.

**Appendix A. Proof of Lemma 1.** To prove the claim in Lemma 1, we first prove the the following lemma which is based on the result in [4, Lemma 5.2].

LEMMA 8. Consider the proposed IQN method in (9). Let  $\mathbf{M}$  be a nonsingular symmetric matrix such that

$$(38) \quad \|\mathbf{M}\mathbf{y}_i^t - \mathbf{M}^{-1}\mathbf{s}_i^t\| \leq \beta \|\mathbf{M}^{-1}\mathbf{s}_i^t\|,$$

for some  $\beta \in [0, 1/3]$  and vectors  $\mathbf{s}_i^t$  and  $\mathbf{y}_i^t$  in  $\mathbb{R}^p$  with  $\mathbf{s}_i^t \neq \mathbf{0}$ . Consider  $i$  as the index of the updated function at step  $t$ , i.e.,  $i = i_t$ , and let  $\mathbf{B}_i^t$  be symmetric and computed according to the update in (15). Then, there exist positive constants  $\alpha$ ,  $\alpha_1$ , and  $\alpha_2$  such that, for any symmetric  $\mathbf{A} \in \mathbb{R}^{p \times p}$  we have,

$$(39) \quad \|\mathbf{B}_i^{t+n} - \mathbf{A}\|_{\mathbf{M}} \leq \left[ (1 - \alpha\theta^2)^{1/2} + \alpha_1 \frac{\|\mathbf{M}\mathbf{y}_i^t - \mathbf{M}^{-1}\mathbf{s}_i^t\|}{\|\mathbf{M}^{-1}\mathbf{s}_i^t\|} \right] \|\mathbf{B}_i^t - \mathbf{A}\|_{\mathbf{M}} + \alpha_2 \frac{\|\mathbf{y}_i^t - \mathbf{A}\mathbf{s}_i^t\|}{\|\mathbf{M}^{-1}\mathbf{s}_i^t\|},$$

where  $\alpha = (1 - 2\beta)/(1 - \beta^2) \in [3/8, 1]$ ,  $\alpha_1 = 2.5(1 - \beta)^{-1}$ ,  $\alpha_2 = 2(1 + 2\sqrt{\beta})\|\mathbf{M}\|_{\mathbf{F}}$ , and

$$(40) \quad \theta = \frac{\|\mathbf{M}(\mathbf{B}_i^t - \mathbf{A})\mathbf{s}_i^t\|}{\|\mathbf{B}_i^t - \mathbf{A}\|_{\mathbf{M}}\|\mathbf{M}^{-1}\mathbf{s}_i^t\|} \quad \text{for } \mathbf{B}_i^t \neq \mathbf{A}, \quad \theta = 0 \quad \text{for } \mathbf{B}_i^t = \mathbf{A}.$$

*Proof.* First note that the Hessian approximation  $\mathbf{B}_i^{t+n}$  is equal to  $\mathbf{B}_i^{t+1}$  if the function  $f_i$  is updated at step  $t$ . Considering this observation and the result of Lemma 5.2. in [4] the claim in (39) follows.  $\square$

The result in Lemma 8 provides an upper bound for the difference between the Hessian approximation matrix  $\mathbf{B}_i^{t+n}$  and any positive definite matrix  $\mathbf{A}$  with respect

to the difference between the previous Hessian approximation  $\mathbf{B}_i^t$  and the matrix  $\mathbf{A}$ . The interesting choice for the arbitrary matrix  $\mathbf{A}$  is the Hessian of the  $i$ -th function at the optimal argument, i.e.,  $\mathbf{A} = \nabla^2 f_i(\mathbf{x}^*)$ , which allows us to capture the difference between the sequence of Hessian approximation matrices for function  $f_i$  and the Hessian  $\nabla^2 f_i(\mathbf{x}^*)$  at the optimal argument. We proceed to use the result in Lemma 8 for  $\mathbf{M} = \nabla^2 f_i(\mathbf{x}^*)^{-1/2}$  and  $\mathbf{A} = \nabla^2 f_i(\mathbf{x}^*)$  to prove the claim in (26). To do so, we first need to show that the condition in (38) is satisfied. Note that according to the condition in Assumptions 1 and 2 we can write

$$(41) \quad \frac{\|\mathbf{y}_i^t - \nabla^2 f_i(\mathbf{x}^*) \mathbf{s}_i^t\|}{\|\nabla^2 f_i(\mathbf{x}^*)^{1/2} \mathbf{s}_i^t\|} \leq \frac{\tilde{L} \|\mathbf{s}_i^t\| \max\{\|\mathbf{z}_i^t - \mathbf{x}^*\|, \|\mathbf{z}_i^{t+1} - \mathbf{x}^*\|\}}{\sqrt{m} \|\mathbf{s}_i^t\|} = \frac{\tilde{L}}{\sqrt{m}} \sigma_i^t$$

This observation implies that the left hand side of the condition in (38) for  $\mathbf{M} = \nabla^2 f_i(\mathbf{x}^*)^{-1/2}$  is bounded above by

$$(42) \quad \frac{\|\mathbf{M} \mathbf{y}_i^t - \mathbf{M}^{-1} \mathbf{s}_i^t\|}{\|\mathbf{M}^{-1} \mathbf{s}_i^t\|} \leq \frac{\|\nabla^2 f_i(\mathbf{x}^*)^{-1/2}\| \|\mathbf{y}_i^t - \nabla^2 f_i(\mathbf{x}^*) \mathbf{s}_i^t\|}{\|\nabla^2 f_i(\mathbf{x}^*)^{1/2} \mathbf{s}_i^t\|} \leq \frac{\tilde{L}}{m} \sigma_i^t$$

Thus, the condition in (38) is satisfied since  $\tilde{L} \sigma_i^t / m < 1/3$ . Replacing the upper bounds in (41) and (42) into the expression in (39) implies the claim in (26) with

$$(43) \quad \beta = \frac{\tilde{L}}{m} \sigma_i^t, \quad \alpha = \frac{1 - 2\beta}{1 - \beta^2}, \quad \alpha_3 = \frac{5\tilde{L}}{2m(1 - \beta)}, \quad \alpha_4 = \frac{2(1 + 2\sqrt{\beta})\tilde{L}}{\sqrt{m}} \|\nabla^2 f_i(\mathbf{x}^*)^{-\frac{1}{2}}\|_{\mathbf{F}},$$

and the proof is complete.

**Appendix B. Proof of Lemma 2.** Start by subtracting  $\mathbf{x}^*$  from both sides of (9) to obtain

$$(44) \quad \mathbf{x}^{t+1} - \mathbf{x}^* = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t \mathbf{z}_i^t - \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{z}_i^t) - \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t \mathbf{x}^* \right).$$

As the gradient of  $f$  at the optimal point is the vector zero, i.e.,  $(1/n) \sum_{i=1}^n \nabla f_i(\mathbf{x}^*) = \mathbf{0}$ , we can subtract  $(1/n) \sum_{i=1}^n \nabla f_i(\mathbf{x}^*)$  from the right hand side of (44) and rearrange terms to obtain

$$(45) \quad \mathbf{x}^{t+1} - \mathbf{x}^* = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t (\mathbf{z}_i^t - \mathbf{x}^*) - \frac{1}{n} \sum_{i=1}^n (\nabla f_i(\mathbf{z}_i^t) - \nabla f_i(\mathbf{x}^*)) \right).$$

The expression in (45) relates the residual at time  $t + 1$  to the previous  $n$  residuals and the Hessian approximations  $\mathbf{B}_i^t$ . To analyze this further, we can replace the Hessian approximations  $\mathbf{B}_i^t$  with the actual Hessians  $\nabla^2 f_i(\mathbf{x}^*)$  and the approximation difference  $\nabla^2 f_i(\mathbf{x}^*) - \mathbf{B}_i^t$ . To do so, we add and subtract  $(1/n) \sum_{i=1}^n \nabla^2 f_i(\mathbf{x}^*) (\mathbf{z}_i^t - \mathbf{x}^*)$  to the right hand side of (45) and rearrange terms to obtain

$$(46) \quad \begin{aligned} \mathbf{x}^{t+1} - \mathbf{x}^* &= \left( \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n [\nabla^2 f_i(\mathbf{x}^*) (\mathbf{z}_i^t - \mathbf{x}^*) - (\nabla f_i(\mathbf{z}_i^t) - \nabla f_i(\mathbf{x}^*))] \right) \\ &\quad + \left( \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n [\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)] (\mathbf{z}_i^t - \mathbf{x}^*) \right). \end{aligned}$$

We proceed to take the norms of both sides and use the triangle inequality to obtain an upper bound on the norm of the residual  $\|\mathbf{x}^{t+1} - \mathbf{x}^*\|$ ,

$$(47) \quad \begin{aligned} \|\mathbf{x}^{t+1} - \mathbf{x}^*\| \leq & \left\| \left( \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t \right)^{-1} \right\| \frac{1}{n} \sum_{i=1}^n \|\nabla^2 f_i(\mathbf{x}^*) (\mathbf{z}_i^t - \mathbf{x}^*) - (\nabla f_i(\mathbf{z}_i^t) - \nabla f_i(\mathbf{x}^*))\| \\ & + \left\| \left( \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t \right)^{-1} \right\| \frac{1}{n} \sum_{i=1}^n \|[\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)] (\mathbf{z}_i^t - \mathbf{x}^*)\|. \end{aligned}$$

To obtain the quadratic term in (30) from the first term in (47), we use the Lipschitz continuity of the Hessians  $\nabla^2 f_i$  which leads to the inequality

$$(48) \quad \|\nabla^2 f_i(\mathbf{x}^*) (\mathbf{z}_i^t - \mathbf{x}^*) - (\nabla f_i(\mathbf{z}_i^t) - \nabla f_i(\mathbf{x}^*))\| \leq \tilde{L} \|\mathbf{z}_i^t - \mathbf{x}^*\|^2.$$

Replacing the expression  $\|\nabla^2 f_i(\mathbf{x}^*) (\mathbf{z}_i^t - \mathbf{x}^*) - (\nabla f_i(\mathbf{z}_i^t) - \nabla f_i(\mathbf{x}^*))\|$  in (47) by the upper bound in (48), the claim in (30) follows.

**Appendix C. Proof of Lemma 3.** In this proof we use some steps in the proof of [4, Theorem 3.2]. To start we use the fact that in a finite-dimensional vector space there always exists a constant  $\eta > 0$  such that  $\|\mathbf{A}\| \leq \eta \|\mathbf{A}\|_{\mathbf{M}}$ . Consider  $\gamma = 1/m$  is an upper bound for the norm  $\|\nabla^2 f(\mathbf{x}^*)^{-1}\|$ . Assume that  $\epsilon(r) = \epsilon$  and  $\delta(r) = \delta$  are chosen such that

$$(49) \quad (2\alpha_3\delta + \alpha_4) \frac{\epsilon}{1-r} \leq \delta \quad \text{and} \quad \gamma(1+r)[\tilde{L}\epsilon + 2\eta\delta] \leq r.$$

Based on the assumption that  $\|\mathbf{B}_i^0 - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} \leq \delta$  we can derive the upper bound  $\|\mathbf{B}_i^0 - \nabla^2 f_i(\mathbf{x}^*)\| \leq \eta\delta$ . This observation along with the inequality  $\|\nabla^2 f_i(\mathbf{x}^*)\| \leq L$  implies that  $\|\mathbf{B}_i^0\| \leq \eta\delta + L$ . Therefore, we obtain  $\|(1/n) \sum_{i=1}^n \mathbf{B}_i^0\| \leq \eta\delta + L$ . The second inequality in (49) implies that  $2\gamma(1+r)\eta\delta \leq r$ . Based on this observation and the inequalities  $\|\mathbf{B}_i^0 - \nabla^2 f_i(\mathbf{x}^*)\| \leq \eta\delta < 2\eta\delta$  and  $\gamma \geq \|\nabla^2 f_i(\mathbf{x}^*)^{-1}\|$ , we obtain from Banach Lemma that  $\|(\mathbf{B}_i^0)^{-1}\| \leq (1+r)\gamma$ . Following the same argument for the matrix  $((1/n) \sum_{i=1}^n \mathbf{B}_i^0)^{-1}$  with the inequalities  $\|(1/n) \sum_{i=1}^n \mathbf{B}_i^0 - (1/n) \sum_{i=1}^n \nabla^2 f_i(\mathbf{x}^*)\| \leq (1/n) \sum_{i=1}^n \|\mathbf{B}_i^0 - \nabla^2 f_i(\mathbf{x}^*)\| \leq \eta\delta$  and  $\|\nabla^2 f(\mathbf{x}^*)^{-1}\| \leq \gamma$  we obtain that

$$(50) \quad \left\| \left( \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^0 \right)^{-1} \right\| \leq (1+r)\gamma.$$

This upper bound in conjunction with the result in (30) yields

$$(51) \quad \begin{aligned} \|\mathbf{x}^1 - \mathbf{x}^*\| & \leq (1+r)\gamma \left[ \frac{\tilde{L}}{n} \sum_{i=1}^n \|\mathbf{z}_i^0 - \mathbf{x}^*\|^2 + \frac{1}{n} \sum_{i=1}^n \|[\mathbf{B}_i^0 - \nabla^2 f_i(\mathbf{x}^*)] (\mathbf{z}_i^0 - \mathbf{x}^*)\| \right] \\ & = (1+r)\gamma \left[ \tilde{L} \|\mathbf{x}^0 - \mathbf{x}^*\|^2 + \frac{1}{n} \sum_{i=1}^n \|[\mathbf{B}_i^0 - \nabla^2 f_i(\mathbf{x}^*)] (\mathbf{x}^0 - \mathbf{x}^*)\| \right]. \end{aligned}$$

Considering the assumptions that  $\|\mathbf{x}^0 - \mathbf{x}^*\| \leq \epsilon$  and  $\|\mathbf{B}_i^0 - \nabla^2 f_i(\mathbf{x}^*)\| \leq \eta\delta < 2\eta\delta$  we can write

$$(52) \quad \begin{aligned} \|\mathbf{x}^1 - \mathbf{x}^*\| & \leq (1+r)\gamma[\tilde{L}\epsilon + 2\eta\delta] \|\mathbf{x}^0 - \mathbf{x}^*\| \\ & \leq r \|\mathbf{x}^0 - \mathbf{x}^*\|, \end{aligned}$$

where the second inequality follows from the second condition in (49). Without loss of generality, assume that  $i_0 = 1$ . Then, based on the result in (26) we obtain

$$\begin{aligned}
\|\mathbf{B}_1^1 - \nabla^2 f_1(\mathbf{x}^*)\|_{\mathbf{M}} &\leq \left[ (1 - \alpha\theta_1^2)^{1/2} + \alpha_3\sigma_1^0 \right] \|\mathbf{B}_1^0 - \nabla^2 f_1(\mathbf{x}^*)\|_{\mathbf{M}} + \alpha_4\sigma_1^0 \\
&\leq (1 + \alpha_3\epsilon)\delta + \alpha_4\epsilon \\
(53) \qquad \qquad \qquad &\leq \delta + 2\alpha_3\epsilon\delta + \alpha_4\epsilon \leq 2\delta.
\end{aligned}$$

We proceed to the next iteration which leads to the inequality

$$\begin{aligned}
\|\mathbf{x}^2 - \mathbf{x}^*\| &\leq (1+r)\gamma \left[ \frac{\tilde{L}}{n} \sum_{i=1}^n \|\mathbf{z}_i^1 - \mathbf{x}^*\|^2 + \frac{1}{n} \sum_{i=1}^n \left\| \left[ \mathbf{B}_i^1 - \nabla^2 f_i(\mathbf{x}^*) \right] (\mathbf{z}_i^1 - \mathbf{x}^*) \right\| \right] \\
&\leq (1+r)\gamma \left[ \tilde{L}\epsilon + 2\eta\delta \right] \left( \frac{n-1}{n} \|\mathbf{x}^0 - \mathbf{x}^*\| + \frac{1}{n} \|\mathbf{x}^1 - \mathbf{x}^*\| \right) \\
&\leq r \left( \frac{n-1}{n} \|\mathbf{x}^0 - \mathbf{x}^*\| + \frac{1}{n} \|\mathbf{x}^1 - \mathbf{x}^*\| \right) \\
(54) \qquad \qquad \qquad &\leq r \|\mathbf{x}^0 - \mathbf{x}^*\|.
\end{aligned}$$

And since the updated index is  $i_1 = 2$  we obtain

$$\begin{aligned}
\|\mathbf{B}_2^2 - \nabla^2 f_2(\mathbf{x}^*)\|_{\mathbf{M}} &\leq \left[ (1 - \alpha\theta_2^2)^{1/2} + \alpha_3\sigma_2^0 \right] \|\mathbf{B}_2^0 - \nabla^2 f_2(\mathbf{x}^*)\|_{\mathbf{M}} + \alpha_4\sigma_2^0 \\
&\leq (1 + \alpha_3\epsilon)\delta + \alpha_4\epsilon \\
(55) \qquad \qquad \qquad &\leq \delta + 2\alpha_3\epsilon\delta + \alpha_4\epsilon \leq 2\delta.
\end{aligned}$$

With the same argument we can show that all  $\|\mathbf{B}_t^t - \nabla^2 f_t(\mathbf{x}^*)\|_{\mathbf{M}} \leq 2\delta$  and  $\|\mathbf{x}^t - \mathbf{x}^*\| \leq \epsilon$ , for all iterates  $t = 1, \dots, n$ . Moreover, we have  $\|\mathbf{x}^t - \mathbf{x}^*\| \leq r \|\mathbf{x}^0 - \mathbf{x}^*\|$  for  $t = 1, \dots, n$ .

Now we use the results for iterates  $t = 1, \dots, n$  as the base of our induction argument. To be more precise, let's assume that for iterates  $t = jn+1, jn+2, \dots, jn+n$  we know that the residuals are bounded above by  $\|\mathbf{x}^t - \mathbf{x}^*\| \leq r^{j+1} \|\mathbf{x}^0 - \mathbf{x}^*\|$  and the Hessian approximation matrices  $\mathbf{B}_i^t$  satisfy the inequalities  $\|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\| \leq 2\eta\delta$ . Our goal is to show that for iterates  $t = (j+1)n+1, (j+1)n+2, \dots, (j+1)n+n$  the inequalities  $\|\mathbf{x}^t - \mathbf{x}^*\| \leq r^{j+2} \|\mathbf{x}^0 - \mathbf{x}^*\|$  and  $\|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\| \leq 2\eta\delta$  hold.

Based on the inequalities  $\|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\| \leq 2\eta\delta$  and  $\|\nabla^2 f_i(\mathbf{x}^*)^{-1}\| \leq \gamma$  we can show that for all  $t = jn+1, jn+2, \dots, jn+n$  we have

$$(56) \qquad \qquad \qquad \left\| \left( \frac{1}{n} \sum_{i=1}^n \mathbf{B}_i^t \right)^{-1} \right\| \leq (1+r)\gamma.$$

Using (56) and the inequality in (26) for the iterate  $t = (j+1)n+1$ , we obtain

$$\begin{aligned}
\|\mathbf{x}^{(j+1)n+1} - \mathbf{x}^*\| &\leq (1+r)\gamma \frac{\tilde{L}}{n} \sum_{i=1}^n \left\| \mathbf{z}_i^{(j+1)n} - \mathbf{x}^* \right\|^2 \\
(57) \qquad \qquad \qquad &+ (1+r)\gamma \frac{1}{n} \sum_{i=1}^n \left\| \left[ \mathbf{B}_i^{(j+1)n} - \nabla^2 f_i(\mathbf{x}^*) \right] \left( \mathbf{z}_i^{(j+1)n} - \mathbf{x}^* \right) \right\|.
\end{aligned}$$

Since the variables are updated in a cyclic fashion the set of variables  $\{\mathbf{z}_i^{(j+1)n}\}_{i=1}^n$  is equal to the set  $\{\mathbf{x}^{(j+1)n-i}\}_{i=0}^{n-1}$ . By considering this relation and replacing the

norms  $\|\mathbf{B}_i^{(j+1)n} - \nabla^2 f_i(\mathbf{x}^*)\|(\mathbf{z}_i^{(j+1)n} - \mathbf{x}^*)\|$  by their upper bounds  $2\eta\delta\|\mathbf{z}_i^{(j+1)n} - \mathbf{x}^*\|$  we can simplify the right hand side of (57) as

$$(58) \quad \|\mathbf{x}^{(j+1)n+1} - \mathbf{x}^*\| \leq (1+r)\gamma \left[ \frac{\tilde{L}}{n} \sum_{i=1}^n \|\mathbf{x}^{jn+i} - \mathbf{x}^*\|^2 + \frac{2\eta\delta}{n} \sum_{i=1}^n \|\mathbf{x}^{jn+i} - \mathbf{x}^*\| \right].$$

Since  $\|\mathbf{x}^{jn+i} - \mathbf{x}^*\| \leq \epsilon$  for all  $j = 1, \dots, n$ , we obtain

$$(59) \quad \|\mathbf{x}^{(j+1)n+1} - \mathbf{x}^*\| \leq (1+r)\gamma \left[ \tilde{L}\epsilon + 2\eta\delta \right] \left( \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}^{jn+i} - \mathbf{x}^*\| \right).$$

According to the second inequality in (49) and the assumption that for iterates  $t = jn+1, jn+2, \dots, jn+n$  we know that  $\|\mathbf{x}^t - \mathbf{x}^*\| \leq r^{j+1}\|\mathbf{x}^0 - \mathbf{x}^*\|$ , we can replace the right hand side of (59) by the following upper bound

$$(60) \quad \|\mathbf{x}^{(j+1)n+1} - \mathbf{x}^*\| \leq r^{j+2}\|\mathbf{x}^0 - \mathbf{x}^*\|.$$

Now we show that the updated Hessian approximation  $\mathbf{B}_{i_t}^{(j+1)n+1}$  for  $t = (j+1)n+1$  satisfies the inequality  $\|\mathbf{B}_{i_t}^{(j+1)n+1} - \nabla^2 f_{i_t}(\mathbf{x}^*)\|_{\mathbf{M}} \leq 2\delta$ . According to the result in (26), we can write

$$(61) \quad \begin{aligned} & \left\| \mathbf{B}_{i_t}^{(j+1)n+1} - \nabla^2 f_{i_t}(\mathbf{x}^*) \right\|_{\mathbf{M}} - \left\| \mathbf{B}_{i_t}^{jn+1} - \nabla^2 f_{i_t}(\mathbf{x}^*) \right\|_{\mathbf{M}} \\ & \leq \alpha_3 \sigma_{i_t}^{jn+1} \left\| \mathbf{B}_{i_t}^{jn+1} - \nabla^2 f_{i_t}(\mathbf{x}^*) \right\|_{\mathbf{M}} + \alpha_4 \sigma_{i_t}^{jn+1}. \end{aligned}$$

Now observe that  $\sigma_{i_t}^{jn+1} = \max\{\|\mathbf{x}^{(j+1)n+1} - \mathbf{x}^*\|, \|\mathbf{x}^{jn+1} - \mathbf{x}^*\|\}$  is bounded above by  $r^{j+1}\|\mathbf{x}^0 - \mathbf{x}^*\|$ . Applying this substitution into (61) and considering the conditions  $\|\mathbf{B}_{i_t}^{jn+1} - \nabla^2 f_{i_t}(\mathbf{x}^*)\|_{\mathbf{M}} \leq 2\delta$  and  $\|\mathbf{x}^0 - \mathbf{x}^*\| \leq \epsilon$  lead to the inequality

$$(62) \quad \left\| \mathbf{B}_{i_t}^{(j+1)n+1} - \nabla^2 f_{i_t}(\mathbf{x}^*) \right\|_{\mathbf{M}} - \left\| \mathbf{B}_{i_t}^{jn+1} - \nabla^2 f_{i_t}(\mathbf{x}^*) \right\|_{\mathbf{M}} \leq r^{j+1}\epsilon(2\delta\alpha_3 + \alpha_4).$$

By writing the expression in (62) for previous iterations and using a recursive logic we obtain that

$$(63) \quad \left\| \mathbf{B}_{i_t}^{(j+1)n+1} - \nabla^2 f_{i_t}(\mathbf{x}^*) \right\|_{\mathbf{M}} - \left\| \mathbf{B}_{i_t}^0 - \nabla^2 f_{i_t}(\mathbf{x}^*) \right\|_{\mathbf{M}} \leq \epsilon(2\delta\alpha_3 + \alpha_4) \frac{1}{1-r}.$$

Based on the first inequality in (49), the right hand side of (63) is bounded above by  $\delta$ . Moreover, the norm  $\|\mathbf{B}_{i_t}^0 - \nabla^2 f_{i_t}(\mathbf{x}^*)\|_{\mathbf{M}}$  is also upper bounded by  $\delta$ . These two bounds imply that

$$(64) \quad \left\| \mathbf{B}_{i_t}^{(j+1)n+1} - \nabla^2 f_{i_t}(\mathbf{x}^*) \right\|_{\mathbf{M}} \leq 2\delta,$$

and consequently  $\|\mathbf{B}_{i_t}^{(j+1)n+1} - \nabla^2 f_{i_t}(\mathbf{x}^*)\| \leq 2\eta\delta$ . By following the steps from (57) to (64), we can show for all iterates  $t = (j+1)n+1, (j+1)n+2, \dots, (j+1)n+n$  the inequalities  $\|\mathbf{x}^t - \mathbf{x}^*\| \leq r^{j+2}\|\mathbf{x}^0 - \mathbf{x}^*\|$  and  $\|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\| \leq 2\eta\delta$  hold. The induction proof is complete and (31) holds. Moreover, the inequality  $\|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\| \leq 2\eta\delta$  holds for all  $i$  and steps  $t$ . Hence, the norms  $\|\mathbf{B}_i^t\|$  and  $\|(\mathbf{B}_i^t)^{-1}\|$ , and consequently  $\|(1/n) \sum_{i=1}^n \mathbf{B}_i^t\|$  and  $\|((1/n) \sum_{i=1}^n \mathbf{B}_i^t)^{-1}\|$  are uniformly bounded.

**Appendix D. Proof of Proposition 4.** According to the result in Lemma 3, we can show that the sequence of errors  $\sigma_i^t = \max\{\|\mathbf{z}_i^{t+1} - \mathbf{x}^*\|, \|\mathbf{z}_i^t - \mathbf{x}^*\|\}$  is summable

for all  $i$ . To do so, consider the sum of the sequence  $\sigma_i^t$  which is upper bounded by

$$(65) \quad \sum_{t=0}^{\infty} \sigma_i^t = \sum_{t=0}^{\infty} \max\{\|\mathbf{z}_i^{t+1} - \mathbf{x}^*\|, \|\mathbf{z}_i^t - \mathbf{x}^*\|\} \leq \sum_{t=0}^{\infty} \|\mathbf{z}_i^{t+1} - \mathbf{x}^*\| + \sum_{t=0}^{\infty} \|\mathbf{z}_i^t - \mathbf{x}^*\|$$

Note that the last time that the index  $i$  is chosen before time  $t$  should be in the set  $\{t-1, \dots, t-n\}$ . This observation in association with the result in (31) implies that

$$(66) \quad \sum_{t=0}^{\infty} \sigma_i^t \leq 2 \sum_{t=0}^{\infty} r^{\lceil \frac{t-n-1}{n} \rceil + 1} \|\mathbf{x}^0 - \mathbf{x}^*\| = 2 \sum_{t=0}^{\infty} r^{\lceil \frac{t-1}{n} \rceil} \|\mathbf{x}^0 - \mathbf{x}^*\|$$

Simplifying the sum in the right hand side of (66) yields

$$(67) \quad \sum_{t=0}^{\infty} \sigma_i^t \leq \frac{2\|\mathbf{x}^0 - \mathbf{x}^*\|}{r} + 2n\|\mathbf{x}^0 - \mathbf{x}^*\| \sum_{t=0}^{\infty} r^t < \infty.$$

Thus, the sequence  $\sigma_i^t$  is summable for all  $i = 1, \dots, n$ . To complete the proof we use the following result from Lemma 3.3 in [12].

LEMMA 9. *Let  $\{\phi^t\}$  and  $\{\delta^t\}$  be sequences of nonnegative numbers such that*

$$(68) \quad \phi^{t+1} \leq (1 + \delta^t)\phi^t + \delta^t \quad \text{and} \quad \sum_{k=1}^{\infty} \delta^k < \infty.$$

*Then, the sequence  $\{\phi^t\}$  converges.*

Considering the results in Lemmata 1 and 9, and the fact that  $\sigma_i^t$  is summable as shown in (67), we obtain that the sequence  $\|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}}$  for  $\mathbf{M} := \nabla^2 f_i(\mathbf{x}^*)^{-1/2}$  is convergent and the following limit exists

$$(69) \quad \lim_{k \rightarrow \infty} \|\nabla^2 f_i(\mathbf{x}^*)^{-1/2} \mathbf{B}_i^k \nabla^2 f_i(\mathbf{x}^*)^{-1/2} - \mathbf{I}\|_{\mathbf{F}} = l$$

where  $l$  is a nonnegative constant. Moreover, following the proof of Theorem 3.4 in [12] we can show that

$$(70) \quad \alpha(\theta_i^t)^2 \|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} \leq \|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} - \|\mathbf{B}_i^{t+1} - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} + \sigma_i^t (\alpha_3 \|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} + \alpha_4),$$

and, therefore, summing both sides implies,

$$(71) \quad \sum_{t=0}^{\infty} (\theta_i^t)^2 \|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} < \infty$$

Replacing  $\theta_i^t$  in (71) by its definition in (27) results in

$$(72) \quad \sum_{t=0}^{\infty} \frac{\|\mathbf{M}(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*))\mathbf{s}_i^t\|^2}{\|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}} \|\mathbf{M}^{-1}\mathbf{s}_i^t\|^2} < \infty$$

Since the norm  $\|\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)\|_{\mathbf{M}}$  is upper bounded and the eigenvalues of the matrix  $\mathbf{M} = \nabla^2 f_i(\mathbf{x}^*)^{-1/2}$  are uniformly lower and upper bounded, we conclude from the result in (72) that

$$(73) \quad \lim_{t \rightarrow \infty} \frac{\|(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*))\mathbf{s}_i^t\|^2}{\|\mathbf{s}_i^t\|^2} = 0,$$

which yields the claim in (32).

**Appendix E. Proof of Lemma 5.** Consider the sets of variable variations  $\mathcal{S}_1 = \{\mathbf{s}_i^{t+n\tau}\}_{\tau=0}^T$  and  $\mathcal{S}_2 = \{\mathbf{s}_i^{t+n\tau}\}_{\tau=0}^\infty$ . It is trivial to show that  $\mathbf{z}_i^t - \mathbf{x}^*$  is in the span of the set  $\mathcal{S}_2$ , since the sequences of variables  $\mathbf{x}^t$  and  $\mathbf{z}_i^t$  converge to  $\mathbf{x}^*$  and we can write  $\mathbf{x}^* - \mathbf{z}_i^t = \sum_{\tau=0}^\infty \mathbf{s}_i^{t+n\tau}$ . We proceed to show that the vector  $\mathbf{z}_i^t - \mathbf{x}^*$  is also in the span of the set  $\mathcal{S}_1$  when  $T$  is sufficiently large. To do so, we use a contradiction argument. Let's assume that the vector  $\mathbf{z}_i^t - \mathbf{x}^*$  does not lie in the span of the set  $\mathcal{S}_1$ , and, therefore, it can be decomposed as the sum of two non-zero vectors given by

$$(74) \quad \mathbf{z}_i^t - \mathbf{x}^* = \mathbf{v}_\parallel^t + \mathbf{v}_\perp^t,$$

where  $\mathbf{v}_\parallel^t$  lies in the span of  $\mathcal{S}_1$  and  $\mathbf{v}_\perp^t$  is orthogonal to the span of  $\mathcal{S}_1$ . Since we assume that  $\mathbf{z}_i^t - \mathbf{x}^*$  does not lie in the span of  $\mathcal{S}_1$ , we obtain that  $\mathbf{z}_i^{t+nT} - \mathbf{x}^*$  also does not lie in this span, since  $\mathbf{z}_i^{t+nT} - \mathbf{x}^*$  can be written as the sum  $\mathbf{z}_i^{t+nT} - \mathbf{x}^* = \mathbf{z}_i^t - \mathbf{x}^* + \sum_{\tau=0}^T \mathbf{s}_i^{t+n\tau}$ . These observations imply that we can also decompose the vector  $\mathbf{z}_i^{t+nT} - \mathbf{x}^*$  as

$$(75) \quad \mathbf{z}_i^{t+nT} - \mathbf{x}^* = \mathbf{v}_\parallel^{t+nT} + \mathbf{v}_\perp^{t+nT},$$

where  $\mathbf{v}_\parallel^{t+nT}$  lies in the span of  $\mathcal{S}_1$  and  $\mathbf{v}_\perp^{t+nT}$  is orthogonal to the span of  $\mathcal{S}_1$ . Moreover, we obtain that  $\mathbf{v}_\perp^{t+nT}$  is equal to  $\mathbf{v}_\perp^t$ , i.e.,

$$(76) \quad \mathbf{v}_\perp^{t+nT} = \mathbf{v}_\perp^t.$$

This is true since  $\mathbf{z}_i^{t+nT} - \mathbf{x}^*$  can be written as the sum of  $\mathbf{z}_i^t - \mathbf{x}^*$  and a group of vectors that lie in the span of  $\mathcal{S}_1$ . We assume that the norm  $\|\mathbf{v}_\perp^{t+nT}\| = \|\mathbf{v}_\perp^t\| = \epsilon$  where  $\epsilon > 0$  is a strictly positive constant. According to the linear convergence of the sequence  $\|\mathbf{x}^t - \mathbf{x}^*\|$  in Lemma 3 we know that

$$(77) \quad \|\mathbf{z}_i^{t+nT} - \mathbf{x}^*\| \leq r^{\lceil \frac{t+nT-1}{n} \rceil + 1} \|\mathbf{x}^0 - \mathbf{x}^*\| \leq r^T \|\mathbf{x}^0 - \mathbf{x}^*\|$$

If we pick large enough  $T$  such that  $r^T \|\mathbf{x}^0 - \mathbf{x}^*\| < \epsilon$ , then we obtain  $\|\mathbf{z}_i^{t+nT} - \mathbf{x}^*\| < \epsilon$  which contradicts the assumption  $\|\mathbf{v}_\perp^t\| = \epsilon$ . Thus, we obtain that the vector  $\mathbf{z}_i^t - \mathbf{x}^*$  is also in the span of set  $\mathcal{S}_1$ .

Since the vector  $\mathbf{z}_i^t - \mathbf{x}^*$  is in the span of  $\mathcal{S}_1$ , we can write the normalized vector  $(\mathbf{z}_i^t - \mathbf{x}^*)/\|\mathbf{z}_i^t - \mathbf{x}^*\|$  as a linear combination of the set of normalized vectors  $\{\mathbf{s}_i^{t+n\tau}/\|\mathbf{s}_i^{t+n\tau}\|\}_{\tau=0}^T$ . This property allows to write

$$(78) \quad \begin{aligned} \lim_{t \rightarrow \infty} \frac{\|(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*))(\mathbf{z}_i^t - \mathbf{x}^*)\|}{\|\mathbf{z}_i^t - \mathbf{x}^*\|} &= \lim_{t \rightarrow \infty} \left\| (\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)) \frac{(\mathbf{z}_i^t - \mathbf{x}^*)}{\|\mathbf{z}_i^t - \mathbf{x}^*\|} \right\| \\ &= \lim_{t \rightarrow \infty} \left\| (\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)) \sum_{\tau=0}^T a_\tau \frac{\mathbf{s}_i^{t+n\tau}}{\|\mathbf{s}_i^{t+n\tau}\|} \right\|, \end{aligned}$$

where  $a_\tau$  is coefficient of the vector  $\mathbf{s}_i^{t+n\tau}$  when we write  $(\mathbf{z}_i^t - \mathbf{x}^*)/\|\mathbf{z}_i^t - \mathbf{x}^*\|$  as the linear combination of the normalized vectors  $\{\mathbf{s}_i^{t+n\tau}/\|\mathbf{s}_i^{t+n\tau}\|\}_{\tau=0}^T$ . Now since the index of the difference  $\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)$  does not match with the descent directions  $\mathbf{s}_i^{t+n\tau}$ . We add and subtract the term  $\mathbf{B}_i^{t+n\tau}$  to the expression  $\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)$  and use the triangle inequality to write

$$(79) \quad \begin{aligned} &\lim_{t \rightarrow \infty} \frac{\|(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*))(\mathbf{z}_i^t - \mathbf{x}^*)\|}{\|\mathbf{z}_i^t - \mathbf{x}^*\|} \\ &\leq \lim_{t \rightarrow \infty} \left\| \sum_{\tau=0}^T a_\tau \frac{(\mathbf{B}_i^{t+n\tau} - \nabla^2 f_i(\mathbf{x}^*))\mathbf{s}_i^{t+n\tau}}{\|\mathbf{s}_i^{t+n\tau}\|} \right\| + \left\| \sum_{\tau=0}^T a_\tau \frac{(\mathbf{B}_i^t - \mathbf{B}_i^{t+n\tau})\mathbf{s}_i^{t+n\tau}}{\|\mathbf{s}_i^{t+n\tau}\|} \right\|. \end{aligned}$$

We first simplify the first limit in the right hand side of (79). Using the Cauchy-Schwarz inequality and the result in Proposition 4 we can write

$$\begin{aligned} \lim_{t \rightarrow \infty} \left\| \sum_{\tau=0}^T a_\tau \frac{(\mathbf{B}_i^{t+n\tau} - \nabla^2 f_i(\mathbf{x}^*)) \mathbf{s}_i^{t+n\tau}}{\|\mathbf{s}_i^{t+n\tau}\|} \right\| &\leq \lim_{t \rightarrow \infty} \sum_{\tau=0}^T a_\tau \left\| \frac{(\mathbf{B}_i^{t+n\tau} - \nabla^2 f_i(\mathbf{x}^*)) \mathbf{s}_i^{t+n\tau}}{\|\mathbf{s}_i^{t+n\tau}\|} \right\| \\ (80) \qquad \qquad \qquad &= \sum_{\tau=0}^T a_\tau \lim_{t \rightarrow \infty} \left\| \frac{(\mathbf{B}_i^{t+n\tau} - \nabla^2 f_i(\mathbf{x}^*)) \mathbf{s}_i^{t+n\tau}}{\|\mathbf{s}_i^{t+n\tau}\|} \right\| = 0. \end{aligned}$$

Based on the results in (79) and (80), to prove the claim in (33) it remains to show

$$(81) \qquad \lim_{t \rightarrow \infty} \left\| \sum_{\tau=0}^T a_\tau \frac{(\mathbf{B}_i^t - \mathbf{B}_i^{t+n\tau}) \mathbf{s}_i^{t+n\tau}}{\|\mathbf{s}_i^{t+n\tau}\|} \right\| = 0.$$

To reach this goal, we first study the limit of the difference between two consecutive update Hessian approximation matrices  $\lim_{t \rightarrow \infty} \|\mathbf{B}_i^t - \mathbf{B}_i^{t+n}\|$ . Note that if we set  $\mathbf{A} = \mathbf{B}_i^t$  in (39), we obtain that

$$(82) \qquad \|\mathbf{B}_i^{t+n} - \mathbf{B}_i^t\|_{\mathbf{M}} \leq \alpha_2 \frac{\|\mathbf{y}_i^t - \mathbf{B}_i^t \mathbf{s}_i^t\|}{\|\mathbf{M}^{-1} \mathbf{s}_i^t\|}.$$

where  $\mathbf{M} = (\nabla^2 f_i(\mathbf{x}^*))^{-1/2}$ . By adding and subtracting the term  $\nabla^2 f_i(\mathbf{x}^*) \mathbf{s}_i^t$  and using the result in (32), we can show that the difference  $\|\mathbf{B}_i^{t+n} - \mathbf{B}_i^t\|_{\mathbf{M}}$  approaches zero asymptotically. In particular,

$$\begin{aligned} \lim_{t \rightarrow \infty} \|\mathbf{B}_i^{t+n} - \mathbf{B}_i^t\|_{\mathbf{M}} &\leq \alpha_2 \lim_{t \rightarrow \infty} \frac{\|\mathbf{y}_i^t - \mathbf{B}_i^t \mathbf{s}_i^t\|}{\|\mathbf{M}^{-1} \mathbf{s}_i^t\|} \\ (83) \qquad \qquad \qquad &\leq \alpha_2 \lim_{t \rightarrow \infty} \frac{\|\mathbf{y}_i^t - \nabla^2 f_i(\mathbf{x}^*) \mathbf{s}_i^t\|}{\|\mathbf{M}^{-1} \mathbf{s}_i^t\|} + \alpha_2 \lim_{t \rightarrow \infty} \frac{\|(\nabla^2 f_i(\mathbf{x}^*) - \mathbf{B}_i^t) \mathbf{s}_i^t\|}{\|\mathbf{M}^{-1} \mathbf{s}_i^t\|}. \end{aligned}$$

Since  $\|\mathbf{y}_i^t - \nabla^2 f_i(\mathbf{x}^*) \mathbf{s}_i^t\|$  is bounded above by  $\tilde{L} \|\mathbf{s}_i^t\| \max\{\|\mathbf{z}_i^t - \mathbf{x}^*\|, \|\mathbf{z}_i^{t+1} - \mathbf{x}^*\|\}$  and the eigenvalues of the matrix  $\mathbf{M}$  are uniformly bounded we obtain that the first limit in the right hand side of (83) converges to zero. Further, the result in (32) shows that the second limit in the right hand side of (83) also converges to zero. Therefore,

$$(84) \qquad \lim_{t \rightarrow \infty} \|\mathbf{B}_i^{t+n} - \mathbf{B}_i^t\|_{\mathbf{M}} = 0.$$

Following the same argument we can show that for any two consecutive Hessian approximation matrices the difference approaches zero asymptotically. Thus, we obtain

$$\begin{aligned} \lim_{t \rightarrow \infty} \|\mathbf{B}_i^t - \mathbf{B}_i^{t+n\tau}\|_{\mathbf{M}} &\leq \lim_{t \rightarrow \infty} \left\| \sum_{u=0}^{\tau-1} (\mathbf{B}_i^{t+nu} - \mathbf{B}_i^{t+n(u+1)}) \right\|_{\mathbf{M}} \\ (85) \qquad \qquad \qquad &\leq \sum_{u=0}^{\tau-1} \lim_{t \rightarrow \infty} \|\mathbf{B}_i^{t+nu} - \mathbf{B}_i^{t+n(u+1)}\|_{\mathbf{M}} = 0. \end{aligned}$$

Observing the result in (85) we can show that

$$\begin{aligned} \lim_{t \rightarrow \infty} \left\| \sum_{\tau=0}^T a_\tau \frac{(\mathbf{B}_i^t - \mathbf{B}_i^{t+n\tau}) \mathbf{s}_i^{t+n\tau}}{\|\mathbf{s}_i^{t+n\tau}\|} \right\| &\leq \sum_{\tau=0}^T a_\tau \lim_{t \rightarrow \infty} \left\| \frac{(\mathbf{B}_i^t - \mathbf{B}_i^{t+n\tau}) \mathbf{s}_i^{t+n\tau}}{\|\mathbf{s}_i^{t+n\tau}\|} \right\| \\ (86) \qquad \qquad \qquad &\leq \sum_{\tau=0}^T a_\tau \lim_{t \rightarrow \infty} \|\mathbf{B}_i^t - \mathbf{B}_i^{t+n\tau}\| = 0. \end{aligned}$$

Therefore, the result in (81) holds. The claim in (33) follows by combining the results in (79), (80), and (81).

**Appendix F. Proof of Theorem 6.** The result in Lemma 2 implies

$$(87) \quad \|\mathbf{x}^{t+1} - \mathbf{x}^*\| \leq \frac{\tilde{L}\Gamma^t}{n} \sum_{i=1}^n \|\mathbf{z}_i^t - \mathbf{x}^*\|^2 + \frac{\Gamma^t}{n} \sum_{i=1}^n \|(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)) (\mathbf{z}_i^t - \mathbf{x}^*)\|.$$

Divide both sides of (87) by  $(1/n) \sum_{i=1}^n \|\mathbf{z}_i^t - \mathbf{x}^*\|$  to obtain

$$(88) \quad \frac{\|\mathbf{x}^{t+1} - \mathbf{x}^*\|}{\frac{1}{n} \sum_{i=1}^n \|\mathbf{z}_i^t - \mathbf{x}^*\|} \leq \tilde{L}\Gamma^t \sum_{i=1}^n \frac{\|\mathbf{z}_i^t - \mathbf{x}^*\|^2}{\sum_{i=1}^n \|\mathbf{z}_i^t - \mathbf{x}^*\|} + \Gamma^t \sum_{i=1}^n \frac{\|(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)) (\mathbf{z}_i^t - \mathbf{x}^*)\|}{\sum_{i=1}^n \|\mathbf{z}_i^t - \mathbf{x}^*\|}$$

Since the error  $\|\mathbf{z}_i^t - \mathbf{x}^*\|$  is a lower bound for the sum of errors  $\sum_{i=1}^n \|\mathbf{z}_i^t - \mathbf{x}^*\|$ , we can replace  $\|\mathbf{z}_i^t - \mathbf{x}^*\|$  for  $\sum_{i=1}^n \|\mathbf{z}_i^t - \mathbf{x}^*\|$  into (88) which implies

$$(89) \quad \begin{aligned} \frac{\|\mathbf{x}^{t+1} - \mathbf{x}^*\|}{\frac{1}{n} \sum_{i=1}^n \|\mathbf{z}_i^t - \mathbf{x}^*\|} &\leq \tilde{L}\Gamma^t \sum_{i=1}^n \frac{\|\mathbf{z}_i^t - \mathbf{x}^*\|^2}{\|\mathbf{z}_i^t - \mathbf{x}^*\|} + \Gamma^t \sum_{i=1}^n \frac{\|(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)) (\mathbf{z}_i^t - \mathbf{x}^*)\|}{\|\mathbf{z}_i^t - \mathbf{x}^*\|} \\ &= \tilde{L}\Gamma^t \sum_{i=1}^n \|\mathbf{z}_i^t - \mathbf{x}^*\| + \Gamma^t \sum_{i=1}^n \frac{\|(\mathbf{B}_i^t - \nabla^2 f_i(\mathbf{x}^*)) (\mathbf{z}_i^t - \mathbf{x}^*)\|}{\|\mathbf{z}_i^t - \mathbf{x}^*\|}. \end{aligned}$$

Since  $\Gamma^t$  is bounded above, computing the limit of both sides in (89) yields

$$(90) \quad \lim_{t \rightarrow \infty} \frac{\|\mathbf{x}^{t+1} - \mathbf{x}^*\|}{\frac{1}{n} \sum_{i=1}^n \|\mathbf{z}_i^t - \mathbf{x}^*\|} = 0.$$

The result in (90) in association with the simplification for the sum  $\sum_{i=1}^n \|\mathbf{z}_i^t - \mathbf{x}^*\| = \sum_{i=0}^{n-1} \|\mathbf{x}^{t-i} - \mathbf{x}^*\|$  leads to the claim in (34).

**Appendix G. Proof of Theorem 7.** Consider the definition of the sequence  $\tilde{\mathbf{x}}^t = \operatorname{argmax}_{u \in \{tn, \dots, tn+n-1\}} \{\|\mathbf{x}^u - \mathbf{x}^*\|\}$  which is a subsequence of the sequence  $\{\mathbf{x}^t\}_{t=0}^\infty$ . Our goal is to show this subsequence converges superlinearly to  $\mathbf{x}^*$ , i.e.,  $\lim_{t \rightarrow \infty} \frac{\|\tilde{\mathbf{x}}^{t+1} - \mathbf{x}^*\|}{\|\tilde{\mathbf{x}}^t - \mathbf{x}^*\|} = 0$ . To do so, first note that the result in Theorem 6 implies that

$$(91) \quad \lim_{t \rightarrow \infty} \frac{\|\mathbf{x}^t - \mathbf{x}^*\|}{\max\{\|\mathbf{x}^{t-1} - \mathbf{x}^*\|, \dots, \|\mathbf{x}^{t-n} - \mathbf{x}^*\|\}} = 0,$$

which follows from the inequality  $\max\{\|\mathbf{x}^{t-1} - \mathbf{x}^*\|, \dots, \|\mathbf{x}^{t-n} - \mathbf{x}^*\|\} \geq (1/n)(\|\mathbf{x}^{t-1} - \mathbf{x}^*\| + \dots + \|\mathbf{x}^{t-n} - \mathbf{x}^*\|)$ . Based on the limit in (91), there exists a large enough  $t_0$  such that for all  $t \geq t_0$  the following inequality holds,

$$(92) \quad \|\mathbf{x}^t - \mathbf{x}^*\| < \max\{\|\mathbf{x}^{t-1} - \mathbf{x}^*\|, \dots, \|\mathbf{x}^{t-n} - \mathbf{x}^*\|\}.$$

Combining the inequality in (92) with the inequalities  $\|\mathbf{x}^{t-i} - \mathbf{x}^*\| \leq \max\{\|\mathbf{x}^{t-1} - \mathbf{x}^*\|, \dots, \|\mathbf{x}^{t-n} - \mathbf{x}^*\|\}$  for  $i = 1, \dots, n-1$  yields

$$(93) \quad \max\{\|\mathbf{x}^t - \mathbf{x}^*\|, \dots, \|\mathbf{x}^{t-n+1} - \mathbf{x}^*\|\} \leq \max\{\|\mathbf{x}^{t-1} - \mathbf{x}^*\|, \dots, \|\mathbf{x}^{t-n} - \mathbf{x}^*\|\},$$

and consequently we can generalize this result to obtain

$$(94) \quad \max\{\|\mathbf{x}^t - \mathbf{x}^*\|, \dots, \|\mathbf{x}^{t-n+1} - \mathbf{x}^*\|\} \leq \max\{\|\mathbf{x}^{t-\tau} - \mathbf{x}^*\|, \dots, \|\mathbf{x}^{t-\tau-n+1} - \mathbf{x}^*\|\},$$

for any positive integer  $\tau$  such that  $t - \tau \geq t_0$ .

We use the result in (94) to build a superlinearly convergent subsequence of the residuals sequence  $\|\mathbf{x}^t - \mathbf{x}^*\|$ . If we define  $\mathbf{x}^{tn+u_i^*}$  as the iterate that has the largest error among the iterates in the  $t + 1$ -th pass, i.e.,

$$(95) \quad \mathbf{x}^{tn+u_i^*} = \underset{u \in \{tn, \dots, tn+n-1\}}{\operatorname{argmax}} \{ \|\mathbf{x}^u - \mathbf{x}^*\| \},$$

then it follows that  $\tilde{\mathbf{x}}^t = \mathbf{x}^{tn+u_i^*}$ , where  $u_i^* \in \{0, 1, \dots, n-1\}$ . Moreover, we obtain

$$(96) \quad \begin{aligned} \frac{\|\tilde{\mathbf{x}}^{t+1} - \mathbf{x}^*\|}{\|\tilde{\mathbf{x}}^t - \mathbf{x}^*\|} &= \frac{\|\mathbf{x}^{(t+1)n+u_{i+1}^*} - \mathbf{x}^*\|}{\max\{\|\mathbf{x}^{tn} - \mathbf{x}^*\|, \dots, \|\mathbf{x}^{tn+n-1} - \mathbf{x}^*\|\}} \\ &\leq \frac{\|\mathbf{x}^{tn+n+u_{i+1}^*} - \mathbf{x}^*\|}{\max\{\|\mathbf{x}^{tn+u_{i+1}^*-1} - \mathbf{x}^*\|, \dots, \|\mathbf{x}^{tn+n+u_{i+1}^*-1} - \mathbf{x}^*\|\}}. \end{aligned}$$

The equality follows from the definition of the iterate  $\tilde{\mathbf{x}}^t$  and the definition in (95), and the inequality holds because of the result in (94). Considering the result in (91), computing the limit of both sides leads to the conclusion that the sequence  $\|\tilde{\mathbf{x}}^t - \mathbf{x}^*\|$  is superlinearly convergent. In other words, we obtain that the subsequence  $\{\|\mathbf{x}^{tn+u_i^*} - \mathbf{x}^*\|\}_{t=0}^{t=\infty}$  superlinearly converges to zero.

Let's define the sequence  $q^t$  such that  $q^{kn} = \dots = q^{kn+n-1} = \|\tilde{\mathbf{x}}^k - \mathbf{x}^*\|$  for  $k = 0, 1, 2, \dots$ , which means that the value of the sequence  $q^t$  is fixed for each pass and is equal to the max error of the corresponding pass. Therefore, it is trivial to show that  $q^t$  is always larger than or equal to  $\|\mathbf{x}^t - \mathbf{x}^*\|$ , i.e.,  $\|\mathbf{x}^t - \mathbf{x}^*\| \leq q^t$  for all  $t \geq 0$ . Now define the sequence  $\zeta^t$  such that  $\zeta^t = q^t$  for  $t = 0, \dots, n-1$ , and for  $t \geq n$

$$(97) \quad \zeta^{kn+i} = q^{kn-1} \left( \frac{q^{kn+n-1}}{q^{kn-1}} \right)^{\frac{i+1}{n}}, \quad \text{for } i = 0, \dots, n-1, \quad k \geq 1.$$

According to this definition we can verify that  $\zeta^t$  is an upper bound for the sequence  $q^t$  and, consequently, an upper bound for the sequence of errors  $\|\mathbf{x}^t - \mathbf{x}^*\|$ . Based on the definition of the sequence  $\zeta^t$  in (97), the ratio  $\zeta^{t+1}/\zeta^t$  is given by  $(q^{\lfloor \frac{t+1}{n} \rfloor n+n-1} / q^{\lfloor \frac{t+1}{n} \rfloor n-1})^{1/n}$ . This simplification in association with the definitions of the sequences  $\|\mathbf{x}^t - \mathbf{x}^*\|$  and  $\|\tilde{\mathbf{x}}^t - \mathbf{x}^*\|$  implies that

$$(98) \quad \lim_{t \rightarrow \infty} \frac{\zeta^{t+1}}{\zeta^t} = \lim_{t \rightarrow \infty} \left( \frac{q^{\lfloor \frac{t+1}{n} \rfloor n+n-1}}{q^{\lfloor \frac{t+1}{n} \rfloor n-1}} \right)^{\frac{1}{n}} = \lim_{t \rightarrow \infty} \left( \frac{\|\tilde{\mathbf{x}}^{\lfloor \frac{t+1}{n} \rfloor} - \mathbf{x}^*\|}{\|\tilde{\mathbf{x}}^{\lfloor \frac{t+1}{n} \rfloor - 1} - \mathbf{x}^*\|} \right)^{\frac{1}{n}} = 0,$$

which leads to the claim in (35).

## REFERENCES

- [1] D. BLATT, A. O. HERO, AND H. GAUCHMAN, *A convergent incremental gradient method with a constant step size*, SIAM Journal on Optimization, 18 (2007), pp. 29–51.
- [2] L. BOTTOU, *Large-scale machine learning with stochastic gradient descent*, In Proceedings of COMPSTAT'2010, (2010), pp. 177–186.
- [3] L. BOTTOU AND Y. L. CUN, *On-line learning for very large datasets*, in Applied Stochastic Models in Business and Industry, vol. 21, pp. 137–151, 2005.
- [4] C. G. BROYDEN, J. E. D. JR., WANG, AND J. J. MORE, *On the local and superlinear convergence of quasi-newton methods*, IMA J. Appl. Math, 12 (1973), pp. 223–245.
- [5] F. BULLO, J. CORTES, AND S. MARTINEZ, *Distributed control of robotic networks: a mathematical approach to motion coordination algorithms*, Princeton University Press, 2009.

- [6] R. H. BYRD, S. HANSEN, J. NOCEDAL, AND Y. SINGER, *A stochastic quasi-Newton method for large-scale optimization*, SIAM Journal on Optimization, 26 (2016), pp. 1008–1031.
- [7] Y. CAO, W. YU, W. REN, AND G. CHEN, *An overview of recent progress in the study of distributed multi-agent coordination*, IEEE Transactions on Industrial Informatics, 9 (2013), pp. 427–438.
- [8] V. CEVHER, S. BECKER, AND M. SCHMIDT, *Convex optimization for big data: Scalable, randomized, and parallel algorithms for big data analytics*, IEEE Signal Processing Magazine, 31 (2014), pp. 32–43.
- [9] A. DEFAZIO, F. R. BACH, AND S. LACOSTE-JULIEN, *SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives*, in Advances in Neural Information Processing Systems 27, Montreal, Quebec, Canada, 2014, pp. 1646–1654.
- [10] R. M. GOWER, D. GOLDFARB, AND P. RICHTÁRIK, *Stochastic block BFGS: Squeezing more curvature out of data*, arXiv preprint arXiv:1603.09649, (2016).
- [11] M. GÜRBÜZBALABAN, A. OZDAGLAR, AND P. PARRILO, *On the convergence rate of incremental aggregated gradient algorithms*, arXiv preprint arXiv:1506.02081, (2015).
- [12] J. J. E. DENNIS AND J. J. MORE, *A characterization of superlinear convergence and its application to quasi-Newton methods*, Mathematics of computation, 28 (1974), pp. 549–560.
- [13] R. JOHNSON AND T. ZHANG, *Accelerating stochastic gradient descent using predictive variance reduction*, in Advances in Neural Information Processing Systems 26, Lake Tahoe, Nevada, United States, 2013, pp. 315–323.
- [14] Y. LECUN, C. CORTES, AND C. J. BURGES, *MNIST handwritten digit database*, AT&T Labs [Online]. Available: <http://yann.lecun.com/exdb/mnist>, (2010).
- [15] C. G. LOPES AND A. H. SAYED, *Diffusion least-mean squares over adaptive networks: Formulation and performance analysis*, IEEE Transactions on Signal Processing, 56 (2008), pp. 3122–3136.
- [16] A. LUCCHI, B. MCWILLIAMS, AND T. HOFMANN, *A variance reduced stochastic Newton method*, arXiv preprint arXiv:1503.08316, (2015).
- [17] A. MOKHTARI, M. GÜRBÜZBALABAN, AND A. RIBEIRO, *Surpassing gradient descent provably: A cyclic incremental method with linear convergence rate*, arXiv preprint arXiv:1611.00347, (2016).
- [18] A. MOKHTARI AND A. RIBEIRO, *RES: Regularized stochastic BFGS algorithm*, IEEE Transactions on Signal Processing, 62 (2014), pp. 6089–6104.
- [19] A. MOKHTARI AND A. RIBEIRO, *Global convergence of online limited memory BFGS*, Journal of Machine Learning Research, 16 (2015), pp. 3151–3181.
- [20] 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, AISTATS 2016, Cadiz, Spain, May 9-11, 2016, 2016, pp. 249–258.
- [21] Y. NESTEROV, *Introductory lectures on convex optimization: A basic course*, vol. 87, Springer Science & Business Media, 2013.
- [22] M. J. D. POWELL, *Some global convergence properties of a variable metric algorithm for minimization without exact line search*, Academic Press, London, UK, 2 ed., 1971.
- [23] A. RIBEIRO, *Ergodic stochastic optimization algorithms for wireless communication and networking*, IEEE Transactions on Signal Processing, 58 (2010), pp. 6369–6386.
- [24] A. RIBEIRO, *Optimal resource allocation in wireless communication and networking*, EURASIP Journal on Wireless Communications and Networking, 2012 (2012), pp. 1–19.
- [25] A. RODOMANOV AND D. KROPOTOV, *A superlinearly-convergent proximal newton-type method for the optimization of finite sums*, in Proceedings of The 33rd International Conference on Machine Learning, 2016, pp. 2597–2605.
- [26] N. L. ROUX, M. W. SCHMIDT, AND F. R. BACH, *A stochastic gradient method with an exponential convergence rate for finite training sets*, in Advances in Neural Information Processing Systems 25, Lake Tahoe, Nevada, United States., 2012, pp. 2672–2680.
- [27] I. SCHIZAS, A. RIBEIRO, AND G. GIANNAKIS, *Consensus in ad hoc wsns with noisy links - part i: Distributed estimation of deterministic signals*, IEEE Transactions on Signal Processing, 56 (2008), pp. 350–364.
- [28] M. SCHMIDT, N. LE ROUX, AND F. BACH, *Minimizing finite sums with the stochastic average gradient*, Mathematical Programming, (2016), pp. 1–30.
- [29] 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, AISTATS 2007, pp. 436–443.
- [30] S. SHALEV-SHWARTZ AND N. SREBRO, *SVM optimization: inverse dependence on training set size*, in Machine Learning, Proceedings of the Twenty-Fifth International Conference (ICML 2008), Helsinki, Finland, 2008, pp. 928–935.