

A Decentralized Second-Order Method with Exact Linear Convergence Rate for Consensus Optimization

Aryan Mokhtari, Wei Shi, Qing Ling, and Alejandro Ribeiro

Abstract—This paper considers decentralized consensus optimization problems where different summands of a global objective function are available at nodes of a network that can communicate with neighbors only. The proximal method of multipliers is considered as a powerful tool that relies on proximal primal descent and dual ascent updates on a suitably defined augmented Lagrangian. The structure of the augmented Lagrangian makes this problem nondecomposable, which precludes distributed implementations. This problem is regularly addressed by the use of the alternating direction method of multipliers. The exact second-order method (ESOM) is introduced here as an alternative that relies on: First, the use of a separable quadratic approximation of the augmented Lagrangian, and second, a truncated Taylor's series to estimate the solution of the first-order condition imposed on the minimization of the quadratic approximation of the augmented Lagrangian. The sequences of primal and dual variables generated by ESOM are shown to converge linearly to their optimal arguments when the aggregate cost function is strongly convex and its gradients are Lipschitz continuous. Numerical results demonstrate advantages of ESOM relative to decentralized alternatives in solving least-squares and logistic regression problems.

Index Terms—Decentralized optimization, method of multipliers, multi-agent networks, second-order methods.

I. INTRODUCTION

IN DECENTRALIZED consensus optimization problems, components of a global objective function that is to be minimized are available at different nodes of a network. Formally, consider a decision variable $\tilde{\mathbf{x}} \in \mathbb{R}^p$ and a connected network containing n nodes where each node i has access to a local objective function $f_i : \mathbb{R}^p \rightarrow \mathbb{R}$. Nodes can exchange information with neighbors only and try to minimize the global

cost function $\sum_{i=1}^n f_i(\tilde{\mathbf{x}})$,

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

We assume that the local objective functions $f_i(\tilde{\mathbf{x}})$ are strongly convex. The global objective function $\sum_{i=1}^n f_i(\tilde{\mathbf{x}})$, which is the sum of a set of strongly convex functions, is also strongly convex. Problems like (1) arise in decentralized control [1]–[3], wireless communication [4], [5], sensor networks [6]–[8], and large scale machine learning [9]–[11].

Decentralized methods for solving (1) can be divided into two classes: primal domain methods and dual domain methods. Decentralized gradient descent (DGD) is a well-established primal method that implements gradient descent on a penalized version of (1) whose gradient can be separated into per-node components. Network Newton (NN) is a more recent alternative that accelerates the convergence of DGD by incorporating second order information of the penalized objective [12], [13]. Both, DGD and NN, converge to a neighborhood of the optimal argument $\tilde{\mathbf{x}}^*$ when using a constant stepsize and converge sublinearly to the exact optimal argument if using a diminishing stepsize.

Dual domain methods build on the fact that the dual function of (1) has a gradient with separable structure. The use of plain dual gradient descent is possible but generally slow to converge [14]–[16]. In centralized optimization, better convergence speeds are attained by the method of multipliers (MM) that adds a quadratic augmentation term to the Lagrangian [17], [18], or the proximal (P)MM that adds an additional term to keep iterates close. In either case, the quadratic term that is added to construct the augmented Lagrangian makes distributed computation of primal gradients impossible. This issue is most often overcome with the use of decentralized (D) versions of the alternating direction method of multipliers (ADMM) [6], [19], [20]. Besides the ADMM, other methods that use different alternatives to approximate the gradients of the dual function have also been proposed [21]–[27]. The convergence rates of these methods have not been studied except for the DADMM and its variants that are known to converge linearly to the optimal argument when the local functions are strongly convex and their gradients are Lipschitz continuous [20], [28], [29]. An important observation here is that while all of these methods try to approximate the MM or the PMM, the performance penalty entailed by the approximation has not been studied.

This paper introduces the exact second order method (ESOM) which uses quadratic approximations of the augmented Lagrangians of (1) and leads to a set of separable subproblems. Similar to other second order methods, implementation

Manuscript received February 1, 2016; revised July 30, 2016; accepted September 19, 2016. Date of publication September 26, 2016; date of current version November 4, 2016. This work was supported in part by the National Science Foundation CAREER CCF-0952867, in part by the ONR under Grant N00014-12-1-0997, in part by the National Science Foundation China under Grant 61573331, and in part by the National Science Foundation Anhui under Grant 1608085QF130. The guest editor coordinating the review of this manuscript and approving it for publication was Prof. Vincenzo Matta.

A. Mokhtari and A. Ribeiro are with the Department of Electrical and Systems Engineering, University of Pennsylvania, Philadelphia, PA 19104 USA (e-mail: aryanm@seas.upenn.edu; aribeiro@seas.upenn.edu).

W. Shi is with the Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, Urbana, IL 61801 USA (e-mail: wilburs@illinois.edu).

Q. Ling is with the Department of Automation, University of Science and Technology of China, Anhui 230026, China (e-mail: qingling@mail.ustc.edu.cn).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSPN.2016.2613678

of ESOM requires computation of Hessian inverses. Distributed implementation of this operation is infeasible because while the Hessian of the proximal augmented Lagrangian is neighbor sparse, its inverse is not. ESOM resolves this issue by using the Hessian inverse approximation technique introduced in [12], [13], [30]. This technique consists of truncating the Taylor's series of the Hessian inverse to order K to obtain the family of methods ESOM- K . Implementation of this expansion in terms of local operations is possible. A remarkable property of all ESOM- K methods is that they can be shown to pay a performance penalty relative to (centralized) PMM that vanishes with increasing iterations.

We begin the paper by reformulating (1) in a form more suitable for decentralized implementation (Proposition 1) and proceed to describe the PMM (Section II). ESOM is a variation of PMM that substitutes the proximal augmented Lagrangian with its quadratic approximation (Section III). Implementation of ESOM requires computing the inverse of the Hessian of the proximal augmented Lagrangian. Since this inversion cannot be computed using local and neighboring information, ESOM- K approximates the Hessian inverse with the K -order truncation of the Taylor's series expansion of the Hessian inverse. This expansion can be carried out using an inner loop of local operations. This and other details required for decentralized implementation of ESOM- K are discussed in Section III-A along with a discussion of how ESOM can be interpreted as a saddle point generalization of the Network Newton methods proposed in [31] (Remark 2) or a second order version of the EXTRA method in [32] (Remark 3).

Convergence analyses of PMM and ESOM are then presented (Section IV). Linear convergence of PMM is established (Section IV-A) and linear convergence factors explicitly derived to use as benchmarks (Theorem 1). In the ESOM analysis (Section IV-B) we provide an upper bound for the error of the proximal augmented Lagrangian approximation (Lemma 3). We leverage this result to prove linear convergence of ESOM (Theorem 2) and to show that ESOM's linear convergence factor approaches the corresponding PMM factor as time grows (Section IV-C). This indicates that the convergence paths of (distributed) ESOM- K and (centralized) PMM are very close. We also study the dependency of the convergence constant with the algorithm's order K .

ESOM tradeoffs and comparisons with other decentralized methods for solving consensus optimization problems are illustrated in numerical experiments (Section V) for a decentralized least squares problem (Section V-A) and a decentralized logistic regression classification problem (Section V-B). Numerical results in both settings verify that larger K leads to faster convergence in terms of number of iterations. However, we observe that all versions of ESOM- K exhibit similar convergence rates in terms of the number of communication exchanges. This implies that ESOM-0 is preferable with respect to the latter metric and that larger K is justified when computational cost is of interest. Faster convergence relative to EXTRA, Network Newton, and DQM is observed. We close the paper with concluding remarks (Section VI).

Notation: Vectors are written as $\mathbf{x} \in \mathbb{R}^n$ and matrices as $\mathbf{A} \in \mathbb{R}^{n \times n}$. Given n vectors \mathbf{x}_i , the vector $\mathbf{x} = [\mathbf{x}_1; \dots; \mathbf{x}_n]$ represents a stacking of the elements of each individual \mathbf{x}_i . We use $\|\mathbf{x}\|$ and $\|\mathbf{A}\|$ to denote the Euclidean norm of vector \mathbf{x} and matrix \mathbf{A} , respectively. The norm of vector \mathbf{x} with respect to positive definite matrix \mathbf{A} is $\|\mathbf{x}\|_{\mathbf{A}} := (\mathbf{x}^T \mathbf{A} \mathbf{x})^{1/2}$. Given a function f its gradient \mathbf{x} is denoted as $\nabla f(\mathbf{x})$ and its Hessian as $\nabla^2 f(\mathbf{x})$.

II. PROXIMAL METHOD OF MULTIPLIERS

Let $\mathbf{x}_i \in \mathbb{R}^p$ be a copy of the decision variable \mathbf{x} kept at node i and define \mathcal{N}_i as the neighborhood of node i . Assuming the network is bidirectionally connected, the optimization problem in (1) is equivalent to the program

$$\begin{aligned} \{\mathbf{x}_i^*\}_{i=1}^n &:= \operatorname{argmin}_{\{\mathbf{x}_i\}_{i=1}^n} \sum_{i=1}^n f_i(\mathbf{x}_i), \\ \text{s.t. } \mathbf{x}_i &= \mathbf{x}_j, \quad \text{for all } i, j \in \mathcal{N}_i. \end{aligned} \quad (2)$$

Indeed, the constraint in (2) enforces the consensus condition $\mathbf{x}_1 = \dots = \mathbf{x}_n$ for any feasible point of (2). With this condition satisfied, the objective in (2) is equal to the objective function in (1) from where it follows that the optimal local variables \mathbf{x}_i^* are all equal to the optimal argument $\tilde{\mathbf{x}}^*$ of (1), i.e., $\mathbf{x}_1^* = \dots = \mathbf{x}_n^* = \tilde{\mathbf{x}}^*$.

To derive ESOM define $\mathbf{x} := [\mathbf{x}_1; \dots; \mathbf{x}_n] \in \mathbb{R}^{np}$ as the concatenation of the local decision variables \mathbf{x}_i and the aggregate function $f : \mathbb{R}^{np} \rightarrow \mathbb{R}$ as $f(\mathbf{x}) = f(\mathbf{x}_1, \dots, \mathbf{x}_n) := \sum_{i=1}^n f_i(\mathbf{x}_i)$ as the sum of all the local functions $f_i(\mathbf{x}_i)$. Introduce the matrix $\mathbf{W} \in \mathbb{R}^{n \times n}$ with elements $w_{ij} \geq 0$ representing a weight that node i assigns to variables of node j . The weight $w_{ij} = 0$ if and only if $j \notin \mathcal{N}_i \cup \{i\}$. The matrix \mathbf{W} is further required to satisfy

$$\mathbf{W}^T = \mathbf{W}, \quad \mathbf{W}\mathbf{1} = \mathbf{1}, \quad \text{null}(\mathbf{I} - \mathbf{W}) = \text{span}(\mathbf{1}). \quad (3)$$

The first condition implies that the weights are symmetric, i.e., $w_{ij} = w_{ji}$. The second condition ensures that the weights of a given node sum up to 1, i.e., $\sum_{j=1}^n w_{ij} = 1$ for all i . Since $\mathbf{W}\mathbf{1} = \mathbf{1}$ we have that $\mathbf{I} - \mathbf{W}$ is rank deficient. The last condition $\text{null}(\mathbf{I} - \mathbf{W}) = \text{span}(\mathbf{1})$ makes the rank of $\mathbf{I} - \mathbf{W}$ exactly equal to $n - 1$ [33].

The matrix \mathbf{W} can be used to reformulate (2) as we show in the following proposition.

Proposition 1: Define the matrix $\mathbf{Z} := \mathbf{W} \otimes \mathbf{I}_p \in \mathbb{R}^{np} \times \mathbb{R}^{np}$ as the Kronecker product of the weight matrix \mathbf{W} and the identity matrix \mathbf{I}_p , and consider the definitions of the global vector $\mathbf{x} := [\mathbf{x}_1; \dots; \mathbf{x}_n]$ and aggregate function $f(\mathbf{x}) := \sum_{i=1}^n f_i(\mathbf{x}_i)$. The optimization problem in (2) is equivalent to

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^{np}} f(\mathbf{x}) \quad \text{s.t. } (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{x} = \mathbf{0}. \quad (4)$$

I.e., $\mathbf{x}^* = [\mathbf{x}_1^*; \dots; \mathbf{x}_n^*]$ with $\{\mathbf{x}_i^*\}_{i=1}^n$ the solution of (2). ■

Proof: We just show that the constraint $((\mathbf{I}_n - \mathbf{W}) \otimes \mathbf{I}_p) \mathbf{x} = (\mathbf{I}_{np} - \mathbf{Z}) \mathbf{x} = \mathbf{0}$ is also a consensus constraint. To do so begin by noticing that since $\mathbf{I} - \mathbf{W}$ is positive semidefinite, $\mathbf{I} - \mathbf{Z} = (\mathbf{I} - \mathbf{W}) \otimes \mathbf{I}_p$ is also positive semidefinite. Therefore, the null space of the square root matrix $(\mathbf{I} - \mathbf{Z})^{1/2}$ is equal

to the null space of $\mathbf{I} - \mathbf{Z}$ and we conclude that satisfying the condition $(\mathbf{I} - \mathbf{Z})^{1/2}\mathbf{x}$ is equivalent to the consensus condition $\mathbf{x}_1 = \dots = \mathbf{x}_n$. This observation in conjunction with the definition of the aggregate function $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x}_i)$ shows that the programs in (4) and (3) are equivalent. In particular, the optimal solution of (4) is $\mathbf{x}^* = [\mathbf{x}_1^*; \dots; \mathbf{x}_n^*]$ with $\{\mathbf{x}_i^*\}_{i=1}^n$ the solution of (2). ■

The formulation in (4) is used to define the proximal method of multipliers (PMM) that we consider in this paper. To do so introduce dual variables $\mathbf{v} \in \mathbb{R}^{np}$ to define the augmented Lagrangian $\mathcal{L}(\mathbf{x}, \mathbf{v})$ of (4) as

$$\mathcal{L}(\mathbf{x}, \mathbf{v}) = f(\mathbf{x}) + \mathbf{v}^T (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{x} + \frac{\alpha}{2} \mathbf{x}^T (\mathbf{I} - \mathbf{Z}) \mathbf{x}, \quad (5)$$

where α is a positive constant. Given the properties of the matrix \mathbf{Z} , the augmentation term $(\alpha/2)\mathbf{x}^T (\mathbf{I} - \mathbf{Z}) \mathbf{x}$ is null when the variable \mathbf{x} is a feasible solution of (4). Otherwise, the inner product is positive and behaves as a penalty for the violation of the consensus constraint.

Introduce a time index $t \in \mathbb{N}$ and define \mathbf{x}_t and \mathbf{v}_t as primal and dual iterates at step t . The primal variable \mathbf{x}_{t+1} is updated by minimizing the sum of the augmented Lagrangian in (5) and the proximal term $(\epsilon/2)\|\mathbf{x} - \mathbf{x}_t\|^2$. We then have that

$$\mathbf{x}_{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^{np}} \left\{ \mathcal{L}(\mathbf{x}, \mathbf{v}_t) + \frac{\epsilon}{2} \|\mathbf{x} - \mathbf{x}_t\|^2 \right\}, \quad (6)$$

where the proximal coefficient $\epsilon > 0$ is a strictly positive constant. The dual variable \mathbf{v}_t is updated by ascending through the gradient of the augmented Lagrangian with respect to the dual variable $\nabla_{\mathbf{v}} \mathcal{L}(\mathbf{x}_{t+1}, \mathbf{v}_t)$ with stepsize α

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \alpha (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{x}_{t+1}. \quad (7)$$

The updates in (6) and (7) for PMM can be considered as a generalization of the method of multipliers (MM), because setting the proximal coefficient $\epsilon = 0$ recovers the updates of MM. The proximal term $(\epsilon/2)\|\mathbf{x} - \mathbf{x}_t\|^2$ is added to keep the updated variable \mathbf{x}_{t+1} close to the previous iterate \mathbf{x}_t . This does not affect convergence guarantees but improves computational stability.

The primal update in (6) may be computationally costly—because it requires solving a convex program—and cannot be implemented in a decentralized manner—because the augmentation term $(1/2\alpha)\mathbf{x}^T (\mathbf{I} - \mathbf{Z}) \mathbf{x}$ in (5) is not separable. In the following section, we propose an approximation of PMM that makes the minimization in (6) computationally economic and separable over nodes of the network. This leads to the set of decentralized updates that define the ESOM algorithm.

III. ESOM: EXACT SECOND-ORDER METHOD

To reduce the computational complexity of (6) and obtain a separable update we introduce a second order approximation of the augmented Lagrangian in (5). Consider then the second order Taylor's expansion $\mathcal{L}(\mathbf{x}, \mathbf{v}_t) \approx \mathcal{L}(\mathbf{x}_t, \mathbf{v}_t) + \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_t, \mathbf{v}_t)^T (\mathbf{x} - \mathbf{x}_t) + (1/2)(\mathbf{x} - \mathbf{x}_t)^T \nabla_{\mathbf{x}}^2 \mathcal{L}(\mathbf{x}_t, \mathbf{v}_t) (\mathbf{x} - \mathbf{x}_t)$ of the augmented Lagrangian with respect to \mathbf{x} centered around $(\mathbf{x}_t, \mathbf{v}_t)$. Using this approximation in lieu of $\mathcal{L}(\mathbf{x}, \mathbf{v}_t)$ in (6)

leads to the primal update

$$\mathbf{x}_{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^{np}} \left\{ \mathcal{L}(\mathbf{x}_t, \mathbf{v}_t) + \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_t, \mathbf{v}_t)^T (\mathbf{x} - \mathbf{x}_t) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_t)^T (\nabla_{\mathbf{x}}^2 \mathcal{L}(\mathbf{x}_t, \mathbf{v}_t) + \epsilon \mathbf{I}) (\mathbf{x} - \mathbf{x}_t) \right\}. \quad (8)$$

The minimization in the right hand side of (8) is of a positive definite quadratic form. Thus, upon defining the Hessian matrix $\mathbf{H}_t \in \mathbb{R}^{np \times np}$ as

$$\mathbf{H}_t := \nabla^2 f(\mathbf{x}_t) + \alpha (\mathbf{I} - \mathbf{Z}) + \epsilon \mathbf{I}, \quad (9)$$

and considering the explicit form of the augmented Lagrangian gradient $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_t, \mathbf{v}_t)$ [cf. (5)] it follows that the variable \mathbf{x}_{t+1} in (8) is given by

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{H}_t^{-1} [\nabla f(\mathbf{x}_t) + (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{v}_t + \alpha (\mathbf{I} - \mathbf{Z}) \mathbf{x}_t]. \quad (10)$$

A fundamental observation here is that the matrix \mathbf{H}_t , which is the Hessian of the objective function in (8), is block neighbor sparse. By block neighbor sparse we mean that the (i, j) th block is non-zero if and only if $j \in \mathcal{N}_i$ or $j = i$. To confirm this claim, observe that $\nabla^2 f(\mathbf{x}_t) \in \mathbb{R}^{np \times np}$ is a block diagonal matrix where its i th diagonal block is the Hessian of the i th local function, $\nabla^2 f_i(\mathbf{x}_{i,t}) \in \mathbb{R}^{p \times p}$. Additionally, matrix $\epsilon \mathbf{I}_{np}$ is a diagonal matrix which implies that the term $\nabla^2 f(\mathbf{x}_t) + \epsilon \mathbf{I}_{np}$ is a block diagonal matrix with blocks $\nabla^2 f_i(\mathbf{x}_{i,t}) + \epsilon \mathbf{I}_p$. Further, it follows from the definition of the matrix \mathbf{Z} that the matrix $\mathbf{I} - \mathbf{Z}$ is neighbor sparse. Therefore, the Hessian \mathbf{H}_t is also neighbor sparse. Although the Hessian \mathbf{H}_t is neighbor sparse, its inverse \mathbf{H}_t^{-1} is not. This observation leads to the conclusion that the update in (10) is not implementable in a decentralized manner, i.e., nodes cannot implement (10) by exchanging information only with their neighbors.

To resolve this issue, we use a Hessian inverse approximation that is built on truncating the Taylor's series of the Hessian inverse \mathbf{H}_t^{-1} as in [12], [30]. To do so, we try to decompose the Hessian as $\mathbf{H}_t = \mathbf{D}_t - \mathbf{B}$ where \mathbf{D}_t is a block diagonal positive definite matrix and \mathbf{B} is a neighbor sparse positive semidefinite matrix. In particular, define \mathbf{D}_t as

$$\mathbf{D}_t := \nabla^2 f(\mathbf{x}_t) + \epsilon \mathbf{I} + 2\alpha (\mathbf{I} - \mathbf{Z}_d), \quad (11)$$

where $\mathbf{Z}_d := \operatorname{diag}(\mathbf{Z})$. Observing the definitions of the matrices \mathbf{H}_t and \mathbf{D}_t and considering the relation $\mathbf{B} = \mathbf{D}_t - \mathbf{H}_t$ we conclude that \mathbf{B} is given by

$$\mathbf{B} := \alpha (\mathbf{I} - 2\mathbf{Z}_d + \mathbf{Z}). \quad (12)$$

Notice that using the decomposition $\mathbf{H}_t = \mathbf{D}_t - \mathbf{B}$ and by factoring $\mathbf{D}_t^{1/2}$, the Hessian inverse can be written as $\mathbf{H}_t^{-1} = \mathbf{D}_t^{-1/2} (\mathbf{I} - \mathbf{D}_t^{-1/2} \mathbf{B} \mathbf{D}_t^{-1/2})^{-1} \mathbf{D}_t^{-1/2}$. Observe that the inverse matrix $(\mathbf{I} - \mathbf{D}_t^{-1/2} \mathbf{B} \mathbf{D}_t^{-1/2})^{-1}$ can be substituted by its Taylor's series $\sum_{u=0}^{\infty} (\mathbf{D}_t^{-1/2} \mathbf{B} \mathbf{D}_t^{-1/2})^u$. Note that this is true if the eigenvalues of the matrix $\mathbf{D}_t^{-1/2} \mathbf{B} \mathbf{D}_t^{-1/2}$ are smaller than 1. We prove in Appendix D that this condition is satisfied. However, computation of the series requires global communication which

is not affordable in decentralized settings. Thus, we approximate the Hessian inverse \mathbf{H}_t^{-1} by truncating the first $K + 1$ terms of its Taylor's series which leads to the Hessian inverse approximation $\tilde{\mathbf{H}}_t^{-1}(K)$,

$$\tilde{\mathbf{H}}_t^{-1}(K) := \mathbf{D}_t^{-1/2} \sum_{u=0}^K \left(\mathbf{D}_t^{-1/2} \mathbf{B} \mathbf{D}_t^{-1/2} \right)^u \mathbf{D}_t^{-1/2}. \quad (13)$$

Notice that the approximate Hessian inverse $\tilde{\mathbf{H}}_t^{-1}(K)$ is K -hop block neighbor sparse, i.e., the (i, j) th block is nonzero if and only if there is at least one path between nodes i and j with length K or smaller.

We introduce the Exact Second-Order Method (ESOM) as a second order method for solving decentralized optimization problems which substitutes the Hessian inverse in update (10) by its K block neighbor sparse approximation $\tilde{\mathbf{H}}_t^{-1}(K)$ defined in (13). Therefore, the primal update of ESOM is

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \tilde{\mathbf{H}}_t^{-1}(K) \left[\nabla f(\mathbf{x}_t) + (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{v}_t + \alpha(\mathbf{I} - \mathbf{Z})\mathbf{x}_t \right]. \quad (14)$$

The ESOM dual update is identical to the update in (7),

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \alpha(\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{x}_{t+1}. \quad (15)$$

Notice that ESOM is different from PMM in approximating the augmented Lagrangian in the primal update of PMM by a second order approximation. Further, ESOM approximates the Hessian inverse of the augmented Lagrangian by truncating the Taylor's series of the Hessian inverse which is not necessarily neighbor sparse. In the following subsection we study the implantation details of the updates in (14) and (15).

Remark 1: The Hessian decomposition $\mathbf{H}_t = \mathbf{D}_t - \mathbf{B}$ with the matrices \mathbf{D}_t and \mathbf{B} in (11) and (12), respectively, is not the only valid decomposition. All decompositions of the form $\mathbf{H}_t = \mathbf{D}_t \pm \mathbf{B}_t$ are valid if \mathbf{D}_t is positive definite and the eigenvalues of the matrix $\mathbf{D}_t^{-1/2} \mathbf{B}_t \mathbf{D}_t^{-1/2}$ are in the interval $(-1, 1)$. The suggested framework guarantees that the matrix \mathbf{B} is positive semidefinite which is helpful in the analysis of the proposed ESOM method. A more comprehensive study of alternative decompositions is studied in [34].

A. Decentralized Implementation of ESOM

The updates in (14) and (15) show that ESOM is a second order approximation of PMM. Although these updates are necessary for understanding the rationale behind ESOM, they are not implementable in a decentralized fashion since the matrix $(\mathbf{I} - \mathbf{Z})^{1/2}$ is not neighbor sparse. To resolve this issue, define the sequence of variables \mathbf{q}_t as $\mathbf{q}_t := (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{v}_t$. Considering the definition of \mathbf{q}_t , the primal update in (14) can be written as

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \tilde{\mathbf{H}}_t^{-1}(K) (\nabla f(\mathbf{x}_t) + \mathbf{q}_t + \alpha(\mathbf{I} - \mathbf{Z})\mathbf{x}_t). \quad (16)$$

Multiplying the dual update in (15) by $(\mathbf{I} - \mathbf{Z})^{1/2}$ from the left hand side and using the definition $\mathbf{q}_t := (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{v}_t$ yields

$$\mathbf{q}_{t+1} = \mathbf{q}_t + \alpha(\mathbf{I} - \mathbf{Z})\mathbf{x}_{t+1}. \quad (17)$$

Notice that the system of updates in (16) and (17) is equivalent to the updates in (14) and (15), i.e., the sequences of variables

\mathbf{x}_t generated by them are identical. Nodes can implement the primal-dual updates in (16) and (17) in a decentralized manner, since the squared root matrix $(\mathbf{I} - \mathbf{Z})^{1/2}$ is eliminated from the updates and nodes can compute the products $(\mathbf{I} - \mathbf{Z})\mathbf{x}_t$ and $(\mathbf{I} - \mathbf{Z})\mathbf{x}_{t+1}$ by exchanging information with their neighbors.

To characterize the local update of each node for implementing the updates in (16) and (17), define

$$\mathbf{g}_t := \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_t, \mathbf{v}_t) = \nabla f(\mathbf{x}_t) + \mathbf{q}_t + \alpha(\mathbf{I} - \mathbf{Z})\mathbf{x}_t, \quad (18)$$

as the gradient of the augmented Lagrangian in (5). Further, define the primal descent direction $\mathbf{d}_t(K)$ with K levels of approximation as

$$\mathbf{d}_t(K) := -\tilde{\mathbf{H}}_t^{-1}(K) \mathbf{g}_t, \quad (19)$$

which implies that the update in (16) can be written as $\mathbf{x}_{t+1} = \mathbf{x}_t + \mathbf{d}_t(K)$. According to the definitions of the Hessian inverse approximation in (13), the explicit expression for the descent direction $\mathbf{d}_t(K)$ is given by $\mathbf{d}_t(K) = \mathbf{D}_t^{-1/2} \sum_{u=0}^K \left(\mathbf{D}_t^{-1/2} \mathbf{B} \mathbf{D}_t^{-1/2} \right)^u \mathbf{D}_t^{-1/2} \mathbf{g}_t$. Considering this definition, we can simplify the expression for the descent direction $\mathbf{d}_t(k+1)$ as

$$\begin{aligned} \mathbf{d}_t(k+1) &= -\mathbf{D}_t^{-1/2} \sum_{u=1}^{k+1} \left(\mathbf{D}_t^{-1/2} \mathbf{B} \mathbf{D}_t^{-1/2} \right)^u \mathbf{D}_t^{-1/2} \mathbf{g}_t \\ &\quad - \mathbf{D}_t^{-1} \mathbf{g}_t, \end{aligned} \quad (20)$$

where we have separated the first term of the sum from the rest. Factorize $\mathbf{D}_t^{-1} \mathbf{B}$ from the summands in (20) to obtain

$$\begin{aligned} \mathbf{d}_t(k+1) &= -\mathbf{D}_t^{-1} \mathbf{B} \mathbf{D}_t^{-1/2} \sum_{u=0}^k \left(\mathbf{D}_t^{-1/2} \mathbf{B} \mathbf{D}_t^{-1/2} \right)^u \mathbf{D}_t^{-1/2} \mathbf{g}_t \\ &\quad - \mathbf{D}_t^{-1} \mathbf{g}_t. \end{aligned} \quad (21)$$

Based on the definition of the descent direction $\mathbf{d}_t(k)$, we obtain that the first term in the right hand side of (21) can be simplified as $\mathbf{D}_t^{-1} \mathbf{B} \mathbf{d}_t(k)$. Therefore, the descent directions $\mathbf{d}_t(k)$ and $\mathbf{d}_t(k+1)$ satisfy the condition

$$\mathbf{d}_t(k+1) = \mathbf{D}_t^{-1} \mathbf{B} \mathbf{d}_t(k) - \mathbf{D}_t^{-1} \mathbf{g}_t. \quad (22)$$

Define $\mathbf{d}_{i,t}(k)$ as the descent direction of node i at step t which is the i th element of the global descent direction $\mathbf{d}_t(k) = [\mathbf{d}_{1,t}(k); \dots; \mathbf{d}_{n,t}(k)]$. Therefore, the localized version of the relation in (22) at node i is given by

$$\mathbf{d}_{i,t}(k+1) = \mathbf{D}_{ii,t}^{-1} \sum_{j=i,j \in \mathcal{N}_i} \mathbf{B}_{ij} \mathbf{d}_{j,t}(k) - \mathbf{D}_{ii,t}^{-1} \mathbf{g}_{i,t}. \quad (23)$$

The update in (23) shows that node i can compute its $(k+1)$ th descent direction $\mathbf{d}_{i,t}(k+1)$ if it has access to the k th descent direction $\mathbf{d}_{i,t}(k)$ of itself and its neighbors $\mathbf{d}_{j,t}(k)$ for $j \in \mathcal{N}_i$. Thus, if nodes initialize with the ESOM-0 descent direction $\mathbf{d}_{i,t}(0) = -\mathbf{D}_{ii,t}^{-1} \mathbf{g}_{i,t}$ and exchange their descent directions with their neighbors for K rounds and use the update in (23), they can compute their local ESOM- K descent direction $\mathbf{d}_{i,t}(K)$. Notice that the i th diagonal block \mathbf{D}_t is given by $\mathbf{D}_{ii,t} := \nabla^2 f_i(\mathbf{x}_{i,t}) + (2\alpha(1 - w_{ii}) + \epsilon)\mathbf{I}$, where $\mathbf{x}_{i,t}$ is the primal variable of node i at step t . Thus, the block $\mathbf{D}_{ii,t}$ is locally available at node i .

Algorithm 1: ESOM- K Method at Node i .

Require: Initial iterates $\mathbf{x}_{i,0} = \mathbf{x}_{j,0} = \mathbf{0}$ for $j \in \mathcal{N}_i$ and $\mathbf{q}_{i,0} = \mathbf{0}$.

- 1: **B** blocks: $\mathbf{B}_{ii} = \alpha(1 - w_{ii})\mathbf{I}$ and $\mathbf{B}_{ij} = \alpha w_{ij}\mathbf{I}$
- 2: **for** $t = 0, 1, 2, \dots$ **do**
- 3: **D** block: $\mathbf{D}_{ii,t} = \nabla^2 f_i(\mathbf{x}_{i,t}) + (2\alpha(1 - w_{ii}) + \epsilon)\mathbf{I}$
- 4: Compute $\mathbf{g}_{i,t} = \nabla f_i(\mathbf{x}_{i,t}) + \mathbf{q}_{i,t} + \alpha(1 - w_{ii})\mathbf{x}_{i,t} - \alpha \sum_{j \in \mathcal{N}_i} w_{ij}\mathbf{x}_{j,t}$
- 5: Compute ESOM-0 descent direction $\mathbf{d}_{i,t}(0) = -\mathbf{D}_{ii,t}^{-1}\mathbf{g}_{i,t}$
- 6: **for** $k = 0, \dots, K - 1$ **do**
- 7: Exchange $\mathbf{d}_{i,t}(k)$ with neighbors $j \in \mathcal{N}_i$
- 8: Compute $\mathbf{d}_{i,t}(k+1) = \mathbf{D}_{ii,t}^{-1} \times [\sum_{j \in \mathcal{N}_i, j=i} \mathbf{B}_{ij}\mathbf{d}_{j,t}(k) - \mathbf{g}_{i,t}]$
- 9: **end for**
- 10: Update primal iterate: $\mathbf{x}_{i,t+1} = \mathbf{x}_{i,t} + \mathbf{d}_{i,t}(K)$.
- 11: Exchange iterates $\mathbf{x}_{i,t+1}$ with neighbors $j \in \mathcal{N}_i$.
- 12: Update dual iterate: $\mathbf{q}_{i,t+1} = \mathbf{q}_{i,t} + \alpha(1 - w_{ii})\mathbf{x}_{i,t+1} - \alpha \sum_{j \in \mathcal{N}_i} w_{ij}\mathbf{x}_{j,t+1}$.
- 13: **end for**

Moreover, node i can evaluate the blocks $\mathbf{B}_{ii} = \alpha(1 - w_{ii})\mathbf{I}$ and $\mathbf{B}_{ij} = \alpha w_{ij}\mathbf{I}$ without extra communication. In addition, nodes can compute the gradient \mathbf{g}_i by communicating with their neighbors. To confirm this claim observe that the i th element of $\mathbf{g}_t = [\mathbf{g}_{1,t}; \dots; \mathbf{g}_{n,t}]$ associated with node i is given by

$$\mathbf{g}_{i,t} := \nabla f_i(\mathbf{x}_{i,t}) + \mathbf{q}_{i,t} + \alpha(1 - w_{ii})\mathbf{x}_{i,t} - \alpha \sum_{j \in \mathcal{N}_i} w_{ij}\mathbf{x}_{j,t}, \quad (24)$$

where $\mathbf{q}_{i,t} \in \mathbb{R}^p$ is the i th element of $\mathbf{q}_t = [\mathbf{q}_{1,t}; \dots; \mathbf{q}_{n,t}]$ and $\mathbf{x}_{i,t}$ the primal variable of node i at step t and they are both available at node i . Hence, the update in (16) can be implemented in a decentralized manner. Likewise, nodes can implement the dual update in (17) using the local update

$$\mathbf{q}_{i,t+1} = \mathbf{q}_{i,t} + \alpha(1 - w_{ii})\mathbf{x}_{i,t+1} - \alpha \sum_{j \in \mathcal{N}_i} w_{ij}\mathbf{x}_{j,t+1}, \quad (25)$$

which requires access to the local primal variable $\mathbf{x}_{j,t+1}$ of the neighboring nodes $j \in \mathcal{N}_i$.

The steps of ESOM- K are summarized in Algorithm 1. The core steps are Steps 5–9 which correspond to computing the ESOM- K primal descent direction $\mathbf{d}_{i,t}(K)$. In Step 5, Each node computes its initial descent direction $\mathbf{d}_{i,t}(0)$ using the block $\mathbf{D}_{ii,t}$ and the local gradient $\mathbf{g}_{i,t}$ computed in Steps 3 and 4, respectively. Steps 7 and 8 correspond to the recursion in (23). In step 7, nodes exchange their k th level descent direction $\mathbf{d}_{i,t}(k)$ with their neighboring nodes to compute the $(k+1)$ th descent direction $\mathbf{d}_{i,t}(k+1)$ in Step 8. The outcome of this recursion is the K th level descent direction $\mathbf{d}_{i,t}(K)$ which is required for the update of the primal variable $\mathbf{x}_{i,t}$ in Step 10. Notice that the blocks of the neighbor sparse matrix \mathbf{B} , which are required for step 8, are computed and stored in Step 1. After updating the primal variables in Step 10, nodes exchange their updated variables $\mathbf{x}_{i,t+1}$ with their neighbors $j \in \mathcal{N}_i$ in Step 11. By

having access to the decision variable of neighboring nodes, nodes update their local dual variable $\mathbf{q}_{i,t}$ in Step 12.

Remark 2: The proposed ESOM algorithm solves problem (4) in the dual domain by defining the proximal augmented Lagrangian. It is also possible to solve problem (4) in the primal domain by solving a penalty version of (4). In particular, by using the quadratic penalty function $(1/2)\|\cdot\|^2$ for the constraint $(\mathbf{I} - \mathbf{Z})^{1/2}\mathbf{x}$ with penalty coefficient α , we obtain the penalized version of (4)

$$\hat{\mathbf{x}}^* := \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^{np}} f(\mathbf{x}) + \frac{\alpha}{2}\mathbf{x}^T(\mathbf{I} - \mathbf{Z})\mathbf{x}, \quad (26)$$

where $\hat{\mathbf{x}}^*$ is the optimal argument of the penalized objective function. Notice that $\hat{\mathbf{x}}^*$ is not equal to the optimal argument \mathbf{x}^* and the distance $\|\mathbf{x}^* - \hat{\mathbf{x}}^*\|$ depends on the choice of α . The objective function in (26) can be minimized by descending through the gradient descent direction which leads to the update of decentralized gradient descent (DGD) [35]. The convergence of DGD can be improved by using Newton's method. Notice that the Hessian of the objective function in (26) is given by

$$\hat{\mathbf{H}} := \nabla^2 f(\mathbf{x}) + \alpha(\mathbf{I} - \mathbf{Z}). \quad (27)$$

The Hessian $\hat{\mathbf{H}}$ in (27) is identical to the Hessian \mathbf{H} in (9) except for the term $\epsilon\mathbf{I}$. Therefore, the same technique for approximating the Hessian inverse $\hat{\mathbf{H}}^{-1}$ can be used to approximate the Newton direction of the penalized objective function in (26) which leads to the update of the Network Newton (NN) methods [12], [13]. Thus, ESOM and NN use an approximate decentralized variation of Newton's method for solving two different problems. In other words, ESOM uses the approximate Newton direction for minimizing the augmented Lagrangian of (4), while NN solves a penalized version of (4) using this approximation. This difference justifies the reason that the sequence of iterates generated by ESOM converges to the optimal argument \mathbf{x}^* (Section IV), while NN converges to a neighborhood of \mathbf{x}^* .

Remark 3: ESOM approximates the augmented Lagrangian $\mathcal{L}(\mathbf{x}, \mathbf{v})$ in (6) by its second order approximation. If we substitute the augmented Lagrangian by its first order approximation we can recover the update of EXTRA proposed in [32]. To be more precise, we can substitute $\mathcal{L}(\mathbf{x}, \mathbf{v}_t)$ in (6) by its first order approximation $\mathcal{L}(\mathbf{x}_t, \mathbf{v}_t) + \nabla \mathcal{L}(\mathbf{x}_t, \mathbf{v}_t)^T(\mathbf{x} - \mathbf{x}_t)$ near the point $(\mathbf{x}_t, \mathbf{v}_t)$ to update the primal variable \mathbf{x} . Considering this substitution, the update of \mathbf{x}_{t+1} is given by

$$\mathbf{x}_{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^{np}} \left\{ \mathcal{L}(\mathbf{x}_t, \mathbf{v}_t) + \nabla \mathcal{L}(\mathbf{x}_t, \mathbf{v}_t)^T(\mathbf{x} - \mathbf{x}_t) + \frac{\epsilon}{2}\|\mathbf{x} - \mathbf{x}_t\|^2 \right\}. \quad (28)$$

Thus, considering the definition of the augmented Lagrangian in (5) the updated variable \mathbf{x}_{t+1} can be explicitly written as

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{1}{\epsilon} \left[\nabla f(\mathbf{x}_t) + (\mathbf{I} - \mathbf{Z})^{1/2}\mathbf{v}_t + \alpha(\mathbf{I} - \mathbf{Z})\mathbf{x}_t \right]. \quad (29)$$

By subtracting the update at step $t-1$ from the update at step t and using the dual variables relation that $\mathbf{v}_{t+1} = \mathbf{v}_t +$

$\alpha(\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{x}_{t+1}$ we obtain the update

$$\begin{aligned} \mathbf{x}_{t+1} = & \left(2\mathbf{I} - \frac{2\alpha}{\epsilon}(\mathbf{I} - \mathbf{Z}) \right) \mathbf{x}_t - \left(\mathbf{I} - \frac{\alpha}{\epsilon}(\mathbf{I} - \mathbf{Z}) \right) \mathbf{x}_{t-1} \\ & - \frac{1}{\epsilon}(\nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}_{t-1})). \end{aligned} \quad (30)$$

The update in (30) shows a first-order approximation of the PMM. It is not hard to show that for specific choices of α and ϵ , the update in (30) is equivalent to the update of EXTRA in [32]. Thus, we expect to observe faster convergence for ESOM relative to EXTRA as it incorporates second-order information. This advantage is studied in Section V.

IV. CONVERGENCE ANALYSIS

In this section, we study convergence rates of PMM and ESOM. First, we show that the sequence of iterates \mathbf{x}_t generated by PMM converges linearly to the optimal argument \mathbf{x}^* . Although, PMM cannot be implemented in a decentralized fashion, its convergence rate can be used as a benchmark for evaluating the performance of ESOM. We then follow the section by analyzing convergence properties ESOM. We show that ESOM exhibits a linear convergence rate and compare its factor of linear convergence with the linear convergence factor of PMM. In proving these results we consider the following assumptions.

Assumption 1: The local objective functions $f_i(\mathbf{x})$ are twice differentiable and the eigenvalues of the local objective functions Hessian $\nabla^2 f(\mathbf{x})$ are bounded by positive constants $0 < m \leq M < \infty$, i.e.

$$m\mathbf{I} \preceq \nabla^2 f_i(\mathbf{x}_i) \preceq M\mathbf{I}, \quad (31)$$

for all $\mathbf{x}_i \in \mathbb{R}^p$ and $i = 1, \dots, n$.

The lower bound in (31) implies that the local objective functions f_i are strongly convex with constant $m > 0$. The upper bound for the eigenvalues of the Hessians $\nabla^2 f_i$ implies that the gradients of the local objective functions ∇f_i are Lipschitz continuous with constant M . Notice that the global objective function $\nabla^2 f(\mathbf{x})$ is a block diagonal matrix where its i th diagonal block is $\nabla^2 f_i(\mathbf{x}_i)$. Therefore, the bounds on the eigenvalues of the local Hessians $\nabla^2 f_i(\mathbf{x}_i)$ in (31) also hold for the global objective function Hessian $\nabla^2 f(\mathbf{x})$. I.e.,

$$m\mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq M\mathbf{I}, \quad (32)$$

for all $\mathbf{x} \in \mathbb{R}^{np}$. Thus, the global objective function f is also strongly convex with constant m and its gradients ∇f are Lipschitz continuous with constant M .

A. Convergence of Proximal Method of Multipliers (PMM)

Convergence rate of PMM can be considered as a benchmark for the convergence rate of ESOM. To establish linear convergence of PMM, We first study the relationship between the primal \mathbf{x} and dual \mathbf{v} iterates generated by PMM and the optimal arguments \mathbf{x}^* and \mathbf{v}^* in the following lemma.

Lemma 1: Consider the updates for the proximal method of multipliers in (6) and (7). The sequences of primal and dual

iterates generated by PMM satisfy

$$\mathbf{v}_{t+1} - \mathbf{v}_t - \alpha(\mathbf{I} - \mathbf{Z})^{1/2}(\mathbf{x}_{t+1} - \mathbf{x}^*) = \mathbf{0}, \quad (33)$$

and

$$\begin{aligned} \nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*) + (\mathbf{I} - \mathbf{Z})^{1/2}(\mathbf{v}_{t+1} - \mathbf{v}^*) \\ + \epsilon(\mathbf{x}_{t+1} - \mathbf{x}_t) = \mathbf{0}. \end{aligned} \quad (34)$$

Proof: See Appendix A. ■

Considering the preliminary results in (33) and (34), we can state convergence results of PMM. To do so, we prove linear convergence of a Lyapunov function of the primal $\|\mathbf{x}_t - \mathbf{x}^*\|^2$ and dual $\|\mathbf{v}_t - \mathbf{v}^*\|^2$ errors. To be more precise, we define the vector $\mathbf{u} \in \mathbb{R}^{2np}$ and matrix $\mathcal{G} \in \mathbb{R}^{np \times np}$ as

$$\mathbf{u} = \begin{bmatrix} \mathbf{v} \\ \mathbf{x} \end{bmatrix}, \quad \mathcal{G} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \alpha\epsilon\mathbf{I} \end{bmatrix}. \quad (35)$$

Notice that the sequence \mathbf{u}_t is the concatenation of the dual variable \mathbf{v}_t and primal variable \mathbf{x}_t . Likewise, we can define \mathbf{u}^* as the concatenation of the optimal arguments \mathbf{v}^* and \mathbf{x}^* . We proceed to prove that the sequence $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2$ converges linearly to null. Observe that $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2$ can be simplified as $\|\mathbf{v}_t - \mathbf{v}^*\|^2 + \alpha\epsilon\|\mathbf{x}_t - \mathbf{x}^*\|^2$. This observation shows that $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2$ is a Lyapunov function of the primal $\|\mathbf{x}_t - \mathbf{x}^*\|^2$ and dual $\|\mathbf{v}_t - \mathbf{v}^*\|^2$ errors. Therefore, linear convergence of the sequence $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2$ implies linear convergence of the sequence $\|\mathbf{x}_t - \mathbf{x}^*\|^2$. In the following theorem, we show that the sequence $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2$ converges to zero at a linear rate.

Theorem 1: Consider the proximal method of multipliers as introduced in (6) and (7). Consider $\beta > 1$ as an arbitrary constant strictly larger than 1 and define $\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})$ as the smallest non-zero eigenvalue of the matrix $\mathbf{I} - \mathbf{Z}$. Further, recall the definitions of the vector \mathbf{u} and matrix \mathcal{G} in (35). If Assumption 1 holds, then the sequence of Lyapunov functions $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2$ generated by PMM satisfies

$$\|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2 \leq \frac{1}{1 + \delta} \|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2, \quad (36)$$

where the constant δ is given by

$$\delta = \min \left\{ \frac{2\alpha\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})}{\beta(m + M)}, \frac{2mM}{\epsilon(m + M)}, \frac{(\beta - 1)\alpha\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})}{\beta\epsilon} \right\}. \quad (37)$$

Proof: See Appendix B. ■

The result in Theorem 1 shows linear convergence of the sequence $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2$ generated by PMM where the factor of linear convergence is $1/(1 + \delta)$. Observe that larger δ implies smaller linear convergence factor $1/(1 + \delta)$ and faster convergence. Notice that all the terms in the minimization in (37) are positive and therefore the constant δ is strictly larger than 0. In addition, the result in Theorem 1 holds for any feasible set of parameters $\beta > 1$, $\epsilon > 0$, and $\alpha > 0$; however, maximizing the parameter δ requires properly choosing the set of parameters β , ϵ , and α .

Observe that when the first positive eigenvalue $\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})$ of the matrix $\mathbf{I} - \mathbf{Z}$, which is the second smallest eigenvalue of

$\mathbf{I} - \mathbf{Z}$, is small the constant δ becomes close to zero and convergence becomes slow. Notice that small $\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})$ shows that the graph is not highly connected. This observation matches the intuition that when the graph has less edges the speed of convergence is slower. Additionally, the upper bounds in (37) show that when the condition number M/m of the global objective function f is large, δ becomes small and the linear convergence becomes slow.

Although PMM enjoys a fast linear convergence rate, each iteration of PMM requires infinite rounds of communications which make it infeasible. In the following section, we study convergence properties of ESOM as a second order approximation of PMM that is implementable in decentralized settings.

B. Convergence of ESOM

We proceed to show that the sequence of iterates \mathbf{x}_t generated by ESOM converges linearly to the optimal argument $\mathbf{x}^* = [\tilde{\mathbf{x}}^*; \dots; \tilde{\mathbf{x}}^*]$. To do so, we first prove linear convergence of the Lyapunov function $\|\mathbf{u}_t - \mathbf{u}^*\|_{\tilde{\mathcal{G}}}^2$ as defined in (35). Moreover, we show that by increasing the Hessian inverse approximation accuracy, ESOM factor of linear convergence can be arbitrary close to the linear convergence factor of PMM in Theorem 1.

Notice that ESOM is built on a second order approximation of the proximal augmented Lagrangian used in the update of PMM. To guarantee that the second order approximation suggested in ESOM is feasible, the local objective functions f_i are required to be twice differentiable as assumed in Assumption 1. The twice differentiability of the local objective functions f_i implies that the aggregate function f , which is the sum of a set of twice differentiable functions, is also twice differentiable. This observation shows that the global objective function $\nabla^2 f(\mathbf{x})$ is definable. Considering this observation, we prove some preliminary results for the iterates generated by ESOM in the following lemma.

Lemma 2: Consider the updates of ESOM in (14) and (15). Recall the definitions of the augmented Lagrangian Hessian $\tilde{\mathbf{H}}_t$ in (9) and the approximate Hessian inverse $\tilde{\mathbf{H}}_t^{-1}(K)$ in (13). If Assumption 1 holds, then the primal and dual iterates generated by ESOM satisfy

$$\mathbf{v}_{t+1} - \mathbf{v}_t - \alpha(\mathbf{I} - \mathbf{Z})^{1/2}(\mathbf{x}_{t+1} - \mathbf{x}^*) = \mathbf{0}. \quad (38)$$

Moreover, we can show that

$$\begin{aligned} \nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*) + (\mathbf{I} - \mathbf{Z})^{1/2}(\mathbf{v}_{t+1} - \mathbf{v}^*) \\ + \epsilon(\mathbf{x}_{t+1} - \mathbf{x}_t) + \mathbf{e}_t = \mathbf{0}, \end{aligned} \quad (39)$$

where the error vector \mathbf{e}_t is defined as

$$\begin{aligned} \mathbf{e}_t := \nabla f(\mathbf{x}_t) + \nabla^2 f(\mathbf{x}_t)(\mathbf{x}_{t+1} - \mathbf{x}_t) - \nabla f(\mathbf{x}_{t+1}) \\ + \left(\tilde{\mathbf{H}}_t(K) - \mathbf{H}_t \right) (\mathbf{x}_{t+1} - \mathbf{x}_t). \end{aligned} \quad (40)$$

Proof: See Appendix C. \blacksquare

The results in Theorem 2 show the relationships between the primal \mathbf{x} and dual \mathbf{v} iterates generated by ESOM and the optimal arguments \mathbf{x}^* and \mathbf{v}^* . The first result in (38) is identical to the convergence property of PMM in (33), while the second result in (39) differs from (34) in having the extra summand \mathbf{e}_t . The

vector \mathbf{e}_t can be interpreted as the error of second order approximation for ESOM at step t . To be more precise, the optimality condition of the primal update of PMM is given by $\nabla f(\mathbf{x}_{t+1}) + (\mathbf{I} - \mathbf{Z})^{1/2}\mathbf{v}_t + \alpha(\mathbf{I} - \mathbf{Z})\mathbf{x}_{t+1} + \epsilon(\mathbf{x}_{t+1} - \mathbf{x}_t) = \mathbf{0}$ as shown in (34). Notice that the second order approximation of this condition is equivalent to $\nabla f(\mathbf{x}_t) + \nabla^2 f(\mathbf{x}_t)(\mathbf{x}_{t+1} - \mathbf{x}_t) + (\mathbf{I} - \mathbf{Z})^{1/2}\mathbf{v}_t + \alpha(\mathbf{I} - \mathbf{Z})\mathbf{x}_{t+1} + \epsilon(\mathbf{x}_{t+1} - \mathbf{x}_t) = \mathbf{0}$. However, the exact Hessian inverse $\mathbf{H}_t^{-1} = (\nabla^2 f(\mathbf{x}_t) + \epsilon\mathbf{I} + \alpha(\mathbf{I} - \tilde{\mathbf{Z}}))^{-1}$ cannot be computed in a distributed manner to solve the optimality condition. Thus, it is approximated by the approximate Hessian inverse matrix $\tilde{\mathbf{H}}_t^{-1}(K)$ as introduced in (13). This shows that the approximate optimality condition in ESOM is $\nabla f(\mathbf{x}_t) + (\mathbf{I} - \mathbf{Z})^{1/2}\mathbf{v}_t + \alpha(\mathbf{I} - \tilde{\mathbf{Z}})\mathbf{x}_t + \tilde{\mathbf{H}}_t(\mathbf{x}_{t+1} - \mathbf{x}_t) = \mathbf{0}$. Hence, the difference between the optimality conditions of PMM and ESOM is $\mathbf{e}_t = \nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}_{t+1}) + \alpha(\mathbf{I} - \tilde{\mathbf{Z}})(\mathbf{x}_t - \mathbf{x}_{t+1}) + \tilde{\mathbf{H}}_t(\mathbf{x}_{t+1} - \mathbf{x}_t) - \epsilon(\mathbf{x}_{t+1} - \mathbf{x}_t)$. By adding and subtracting the term $\mathbf{H}_t(\mathbf{x}_{t+1} - \mathbf{x}_t)$, the definition of the error vector \mathbf{e}_t in (40) follows.

The observation that the vector \mathbf{e}_t characterizes the error of second order approximation in ESOM, motivates analyzing an upper bound for the error vector norm $\|\mathbf{e}_t\|$. To prove that the norm $\|\mathbf{e}_t\|$ is bounded above we assume the following condition is satisfied.

Assumption 2: The global objective function Hessian $\nabla^2 f(\mathbf{x})$ is Lipschitz continuous with constant L , i.e.,

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\tilde{\mathbf{x}})\| \leq L\|\mathbf{x} - \tilde{\mathbf{x}}\|. \quad (41)$$

The conditions imposed by Assumption 2 are customary in the analysis of second-order methods; see, e.g., [29]. In the following lemma, we use the assumption in (41) to prove an upper bound for the error norm $\|\mathbf{e}_t\|$ in terms of $\|\mathbf{x}_{t+1} - \mathbf{x}_t\|$.

Lemma 3: Consider ESOM as introduced in (8) (15) and recall the definition of the error vector \mathbf{e}_t in (40). Further, define $c > 0$ as a lower bound for the local weights w_{ii} . If Assumptions 1 2 hold, then the error vector norm $\|\mathbf{e}_t\|$ is bounded above by

$$\|\mathbf{e}_t\| \leq \Gamma_t \|\mathbf{x}_{t+1} - \mathbf{x}_t\|, \quad (42)$$

where Γ_t is defined as

$$\Gamma_t := \min \left\{ 2M, \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\| \right\} + (M + \epsilon + 2\alpha(1-c))\rho^{K+1}, \quad (43)$$

and $\rho := 2\alpha(1-c)/(2\alpha(1-c) + m + \epsilon)$.

Proof: See Appendix D. \blacksquare

First, note that the lower bound $c > 0$ on the local weights w_{ii} is implied from the fact that all the local weights are positive. In particular, we can define the lower bound c as $c := \min_i w_{ii}$. The result in (42) shows that the error of second order approximation in ESOM vanishes as the sequence of iterates \mathbf{x}_t approaches the optimal argument \mathbf{x}^* . We will show in Theorem 2 that $\|\mathbf{x}_t - \mathbf{x}^*\|$ converges to zero which implies that the limit of the sequence $\|\mathbf{x}_{t+1} - \mathbf{x}_t\|$ is zero.

To understand the definition of Γ_t in (43), we have to decompose the error vector \mathbf{e}_t in (40) into two parts. The first part is $\nabla f(\mathbf{x}_t) + \nabla^2 f(\mathbf{x}_t)(\mathbf{x}_{t+1} - \mathbf{x}_t) - \nabla f(\mathbf{x}_{t+1})$ which comes from the fact that ESOM minimizes a second order approximation of the proximal augmented Lagrangian instead

of the exact proximal augmented Lagrangian. This term can be bounded by $\min\{2M, (L/2)\|\mathbf{x}_{t+1} - \mathbf{x}_t\|\}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|$ as shown in Lemma 3. The second part of the error vector \mathbf{e}_t is $(\tilde{\mathbf{H}}_t(K) - \mathbf{H}_t)(\mathbf{x}_{t+1} - \mathbf{x}_t)$ which shows the error of Hessian inverse approximation. Notice that computation of the exact Hessian inverse \mathbf{H}_t^{-1} is not possible and ESOM approximates the exact Hessian by the approximation $\tilde{\mathbf{H}}_t^{-1}(K)$. According to the results in [12], the difference $\|\tilde{\mathbf{H}}_t(K) - \mathbf{H}_t\|$ can upper bounded by $(M + \epsilon + 2(1 - c)/\alpha)\rho^{K+1}$ which justifies the second term of the expression for Γ_t in (43). In the following theorem, we use the result in Lemma 3 to show that the sequence of Lyapunov functions $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2$ generated by ESOM converges to zero linearly.

Theorem 2: Consider ESOM as introduced in (8) (15). Consider $\beta > 1$ and $\phi > 1$ as arbitrary constants that are strictly larger than 1, and ζ as a positive constant that is chosen from the interval $\zeta \in ((m + M)/2mM, \epsilon/\Gamma_t^2)$. Further, recall the definitions of the vector \mathbf{u} and matrix \mathcal{G} in (35) and consider $\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})$ as the smallest non-zero eigenvalue of the matrix $\mathbf{I} - \mathbf{Z}$. If Assumptions 1 and 2 hold, then the sequence of Lyapunov functions $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2$ generated by ESOM satisfies

$$\|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2 \leq \frac{1}{1 + \delta'_t} \|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2. \quad (44)$$

where the sequence δ'_t is given by

$$\delta'_t = \min \left\{ \frac{2\alpha\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})}{\phi\beta(m + M)}, \left[\frac{2mM}{\epsilon(m + M)} - \frac{1}{\zeta\epsilon} \right], \right. \\ \left. \frac{(\beta - 1)\alpha\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})}{\beta\epsilon} \left[1 - \frac{\zeta\Gamma_t^2}{\epsilon} \right] \left[1 + \frac{\phi\Gamma_t^2(\beta - 1)}{(\phi - 1)\epsilon^2} \right]^{-1} \right\}. \quad (45)$$

Proof: See Appendix E. \blacksquare

The result in Theorem 2 shows linear convergence of the sequence $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2$ generated by ESOM where the factor of linear convergence is $1/(1 + \delta')$. Notice that the positive constant ζ is chosen from the interval $((m + M)/2mM, \epsilon/\Gamma_t^2)$. This interval is non-empty if and only if the proximal parameter ϵ satisfies the condition $\epsilon > \Gamma_t^2(m + M)/2mM$. However, Γ_t also depends on ϵ which makes it unclear if there always exists a choice of ϵ that satisfies the inequality $\epsilon > \Gamma_t^2(m + M)/2mM$. In Appendix F, we provide the condition on ϵ that guarantees $\epsilon > \Gamma_t^2(m + M)/2mM$ holds.

It follows from the result in Theorem 2 that the sequence of primal variables \mathbf{x}_t converges to the optimal argument \mathbf{x}^* defined in (4).

Corollary 1: Under the assumptions in Theorem 2, the sequence of squared errors $\|\mathbf{x}_t - \mathbf{x}^*\|^2$ generated by ESOM converges to zero at a linear rate, i.e.,

$$\|\mathbf{x}_t - \mathbf{x}^*\|^2 \leq \left(\frac{1}{1 + \min_t \{\delta'_t\}} \right)^t \frac{\|\mathbf{u}_0 - \mathbf{u}^*\|_{\mathcal{G}}^2}{\alpha\epsilon}. \quad (46)$$

Proof: According to the definition of the sequence \mathbf{u}_t and matrix \mathcal{G} , we can write $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2 = \alpha\epsilon\|\mathbf{x}_t - \mathbf{x}^*\|^2 + \|\mathbf{v}_t - \mathbf{v}^*\|^2$ which implies that $\|\mathbf{x}_t - \mathbf{x}^*\|^2 \leq (1/\alpha\epsilon)\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2$. Considering this result and linear convergence of the sequence $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2$ in (44), the claim in (46) follows. \blacksquare

C. Convergence Rates Comparison

The expression for δ'_t in (45) verifies the intuition that the convergence rate of ESOM is slower than PMM. This is true, since the upper bounds for δ in PMM are larger than their equivalent upper bounds for δ'_t in ESOM. We obtain that δ'_t is smaller than δ which implies that the linear convergence factor $1/(1 + \delta)$ of PMM is smaller than $1/(1 + \delta'_t)$ for ESOM. Therefore, for all steps t , the linear convergence of PMM is faster than ESOM. Although, linear convergence factor of ESOM $1/(1 + \delta'_t)$ is larger than $1/(1 + \delta)$ for PMM, as time passes the gap between these two constants becomes smaller. In particular, after a number of iterations $(L/2)\|\mathbf{x}_{t+1} - \mathbf{x}_t\|$ becomes smaller than $2M$, and Γ_t can be simplified as

$$\Gamma_t \leq \frac{L}{2}\|\mathbf{x}_{t+1} - \mathbf{x}_t\| + (2\alpha(1 - c) + M + \epsilon)\rho^{K+1}. \quad (47)$$

The term $(L/2)\|\mathbf{x}_{t+1} - \mathbf{x}_t\|$ eventually approaches zero, while the second term $(2(1 - c)/\alpha + M + \epsilon)\rho^{K+1}$ is constant. Although, the second term is not approaching zero, by proper choice of ρ and K , this term can become arbitrary close to zero. Notice that when Γ_t approaches zero, if we set $\zeta = 1/\Gamma_t$ the upper bounds in (45) for δ'_t approach the upper bounds for δ of PMM in (37).

Therefore, as time passes Γ_t becomes smaller, and the factor of linear convergence for ESOM $1/(1 + \delta'_t)$ becomes closer to the linear convergence factor of PMM $1/(1 + \delta)$.

V. NUMERICAL EXPERIMENTS

In this section, we compare the performances of ESOM, EXTRA, Decentralized Quadratically approximated ADMM (DQM), and Network Newton (NN). First, we consider a linear least squares problem and then we use the mentioned methods to solve a logistic regression problem.

A. Decentralized Linear Least Squares

Consider a decentralized linear least squares problem where each agent $i \in \{1, \dots, n\}$ holds its private measurement equation, $\mathbf{y}_i = \mathbf{M}_i\tilde{\mathbf{x}} + \mathbf{v}_i$, where $\mathbf{y}_i \in \mathbb{R}^{m_i}$ and $\mathbf{M}_i \in \mathbb{R}^{m_i \times p}$ are measured data, $\tilde{\mathbf{x}} \in \mathbb{R}^p$ is the unknown variable, and $\mathbf{v}_i \in \mathbb{R}^{m_i}$ is some unknown noise. The decentralized linear least squares estimates $\tilde{\mathbf{x}}$ by solving the optimization problem

$$\tilde{\mathbf{x}}^* = \underset{\tilde{\mathbf{x}}}{\operatorname{argmin}} \sum_{i=1}^n \|\mathbf{M}_i\tilde{\mathbf{x}} - \mathbf{y}_i\|_2^2. \quad (48)$$

The network in this experiment is randomly generated with connectivity ratio $r = 3/n$, where r is defined as the number of edges divided by the number of all possible ones, $n(n - 1)/2$. We set $n = 20$, $p = 5$, and $m_i = 5$ for all $i = 1, \dots, n$. The vectors \mathbf{y}_i and matrices \mathbf{M}_i as well as the noise vectors \mathbf{v}_i , for all i are generated following the standard normal distribution. We precondition the aggregated data matrices \mathbf{M}_i so that the condition number of the problem is 10. The decision variables \mathbf{x}_i are initialized as $\mathbf{x}_{i,0} = 0$ for all nodes $i = 1, \dots, n$ and the initial distance to the optimal is $\|\mathbf{x}_{i,0} - \tilde{\mathbf{x}}^*\| = 100$.

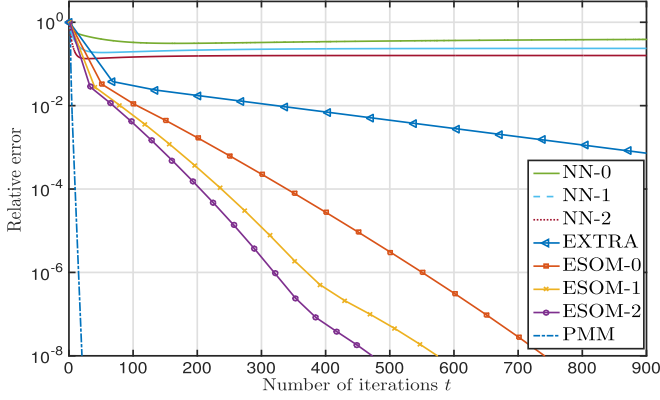


Fig. 1. Relative error $\|\mathbf{x}_t - \mathbf{x}^*\|/\|\mathbf{x}_0 - \mathbf{x}^*\|$ of EXTRA, ESOM- K , NN- K , and PMM versus number of iterations for the least squares problem. Using a larger K for ESOM- K leads to faster convergence and makes the convergence path closer to the one for PMM.

We use Metropolis constant edge weight matrix as the mixing matrix \mathbf{W} in all experiments. We run PMM, EXTRA, and ESOM- K with fixed hand-optimized stepsizes α . The best choices of α for ESOM-0, ESOM-1, and ESOM-2 are $\alpha = 0.03$, $\alpha = 0.04$, and $\alpha = 0.05$, respectively. The stepsize $\alpha = 0.1$ leads to the best performance for EXTRA which is considered in the numerical experiments. Notice that for variations of NN- K , there is no optimal choice of stepsize – smaller stepsize leads to more accurate but slow convergence, while large stepsize accelerates the convergence but to a less accurate neighborhood of the optimal solution. Therefore, for NN-0, NN-1, and NN-2 we set $\alpha = 0.001$, $\alpha = 0.008$, and $\alpha = 0.02$, respectively. Although the PMM algorithm is not implementable in a decentralized fashion, we use its convergence path—which is generated in a centralized manner—as our benchmark. The choice of stepsize for PMM is $\alpha = 2$.

Fig. 1 illustrates the relative error $\|\mathbf{x}_t - \mathbf{x}^*\|/\|\mathbf{x}_0 - \mathbf{x}^*\|$ versus the number of iterations. Notice that the vector \mathbf{x}_t is the concatenation of the local vectors $\mathbf{x}_{i,t}$ and the optimal vector \mathbf{x}^* is defined as $\mathbf{x}^* = [\tilde{\mathbf{x}}^*; \dots; \tilde{\mathbf{x}}^*] \in \mathbb{R}^{np}$. Observe that all the variations of NN- K fail to converge to the optimal argument and they converge linearly to a neighborhood of the optimal solution \mathbf{x}^* . Among the decentralized algorithms with exact linear convergence rate, EXTRA has the worst performance and all the variations of ESOM- K outperform EXTRA. Recall that the problem condition number is 10 in our experiment and the difference between EXTRA and ESOM- K is more significant for problems with larger condition numbers. Further, choosing a larger value of K for ESOM- K leads to faster convergence and as we increase K the convergence path of ESOM- K approaches the convergence path of PMM.

EXTRA requires one round of communications per iteration, while NN- K and ESOM- K require $K + 1$ rounds of local communications per iteration. Thus, convergence paths of these methods in terms of rounds of communications might be different from the ones in Fig. 1. The convergence paths of NN, ESOM, EXTRA in terms of rounds of local communications are shown in Fig. 2. In this plot we ignore PMM, since it requires infinite rounds of communications per iteration. The main

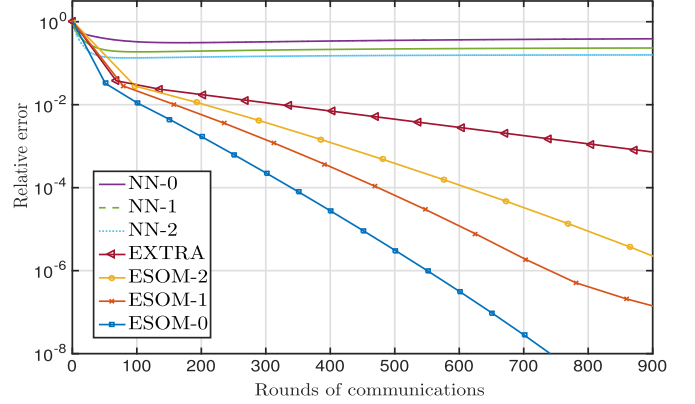


Fig. 2. Relative error $\|\mathbf{x}_t - \mathbf{x}^*\|/\|\mathbf{x}_0 - \mathbf{x}^*\|$ of EXTRA, ESOM- K , NN- K , and PMM versus rounds of communications with neighboring nodes for the least squares problem. ESOM-0 is the most efficient algorithm in terms of communication cost among all the methods.

difference between Figs. 1 and 2 is in the performances of ESOM-0, ESOM-1, and ESOM-2. All of the variations of ESOM outperform EXTRA in terms of rounds of communications, while the best performance belongs to ESOM-0. This observation shows that increasing the approximation level K does not necessary improve the performance of ESOM- K in terms of communication cost.

B. Decentralized Logistic Regression

We consider the application of ESOM for solving a logistic regression problem in a form

$$\tilde{\mathbf{x}}^* := \operatorname{argmin}_{\tilde{\mathbf{x}} \in \mathbb{R}^p} \frac{\lambda}{2} \|\tilde{\mathbf{x}}\|^2 + \sum_{i=1}^n \sum_{j=1}^{m_i} \ln(1 + \exp(-(\mathbf{s}_{ij}^T \tilde{\mathbf{x}}) y_{ij})), \quad (49)$$

where every agent i has access to m_i training samples $(\mathbf{s}_{ij}, y_{ij}) \in \mathbb{R}^p \times \{-1, +1\}$, $j = 1, \dots, m_i$, including explanatory/feature variables \mathbf{s}_{ij} and binary outputs/outcomes y_{ij} . The regularization term $(\lambda/2)\|\tilde{\mathbf{x}}\|^2$ is added to avoid overfitting where $\lambda > 0$. Hence, in the decentralized setting the local objective function f_i of node i is given by

$$f_i(\tilde{\mathbf{x}}) = \frac{\lambda}{2n} \|\tilde{\mathbf{x}}\|^2 + \sum_{j=1}^{m_i} \ln(1 + \exp(-(\mathbf{s}_{ij}^T \tilde{\mathbf{x}}) y_{ij})). \quad (50)$$

The settings are as follows. The connected network is randomly generated with $n = 20$ agents and connectivity ratio $r = 3/n$. Each agent holds 3 samples, i.e., $m_i = 3$, for all i . The dimension of sample vectors \mathbf{s}_{ij} is $p = 3$. The samples are randomly generated, and the optimal logistic classifier $\tilde{\mathbf{x}}^*$ is pre-computed through centralized adaptive gradient method. We use Metropolis constant edge weight matrix as the mixing matrix \mathbf{W} in ESOM- K . The stepsize α for ESOM-0, ESOM-1, ESOM-2, EXTRA, and DQM are hand-optimized and the best of each is used for the comparison.

Figs. 3 and 4 showcase the convergence paths of ESOM-0, ESOM-1, ESOM-2, EXTRA, and DQM versus number of iterations and rounds of communications, respectively. The results match the observations for the least squares problem in

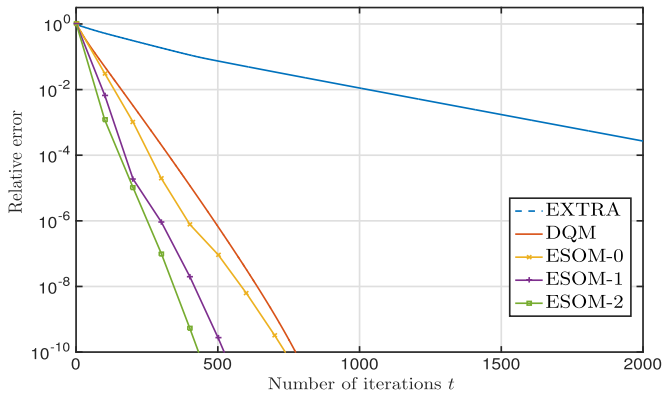


Fig. 3. Relative error $\|\mathbf{x}_t - \mathbf{x}^*\|/\|\mathbf{x}_0 - \mathbf{x}^*\|$ of EXTRA, ESOM- K , and DQM versus number of iterations for the logistic regression problem. EXTRA is significantly slower than the ESOM methods. The proposed methods (ESOM- K) outperform DQM.

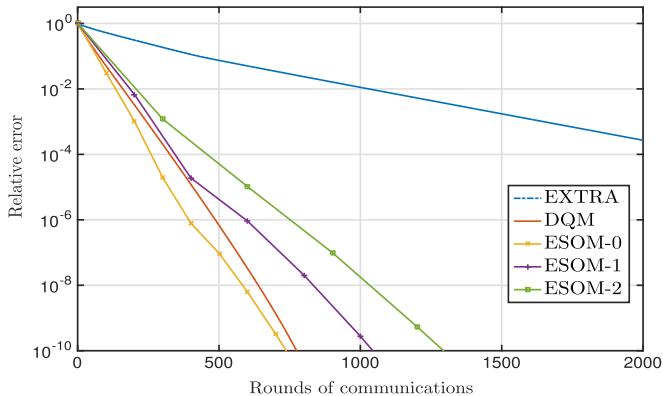


Fig. 4. Relative error $\|\mathbf{x}_t - \mathbf{x}^*\|/\|\mathbf{x}_0 - \mathbf{x}^*\|$ of EXTRA, ESOM- K , and DQM versus rounds of communications for the logistic regression problem. ESOM-0 has the best performance in terms of rounds of communications and it outperforms DQM.

Figs. 1 and 2. Different versions of ESOM- K converge faster than EXTRA both in terms of communication cost and number of iterations. Moreover, ESOM-2 converges faster than ESOM-1 and ESOM-0 in terms of number of iterations, while ESOM-0 has the best performance in terms of communication cost for achieving a target accuracy. Comparing the convergence paths of ESOM-0, ESOM-1, and ESOM-2 with DQM shows that number of iterations required for the convergence of DQM is larger than the required iterations for ESOM-0, ESOM-1, and ESOM-2. In terms of communication cost, DQM has a better performance relative to ESOM-1 and ESOM-2, while ESOM-0 is the most efficient algorithm.

VI. CONCLUSION

We studied the consensus optimization problem where the components of a global objective function are available at different nodes of a network. We proposed an Exact Second-Order Method (ESOM) that converges to the optimal argument of the global objective function at a linear rate. We developed the update of ESOM by substituting the primal update of Proximal Method of Multipliers (PMM) with its second order approximation. Moreover, we approximated the Hessian inverse of the

proximal augmented Lagrangian by truncating its Taylor's series. This approximation leads to a class of algorithms ESOM- K where $K + 1$ indicates the number of Taylor's series terms that are used for Hessian inverse approximation. Convergence analysis of ESOM- K shows that the sequence of iterates converges to the optimal argument linearly irrespective to the choice of K . We showed that the linear convergence factor of ESOM- K is a function of time and the choice of K . The linear convergence factor of ESOM approaches the linear convergence factor of PMM as time passes. Moreover, larger choice of K makes the factor of linear convergence for ESOM closer to the one for PMM. Numerical results verify the theoretical linear convergence and the relation between the linear convergence factor of ESOM- K and PMM. Further, we observed that larger choice of K for ESOM- K leads to faster convergence in terms of number of iterations, while the most efficient version of ESOM- K in terms of communication cost is ESOM-0.

APPENDIX A PROOF OF LEMMA 1

Consider the updates of PMM in (6) and (7). According to (4), the optimal argument \mathbf{x}^* satisfies the condition $(\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{x}^* = \mathbf{0}$. This observation in conjunction with the dual variable update in (7) yields the claim in (33).

To prove the claim in (34), note that the optimality condition of (6) implies $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_{t+1}, \mathbf{v}_t) + \epsilon(\mathbf{x}_{t+1} - \mathbf{x}_t) = \mathbf{0}$. Based on the definition of the Lagrangian $\mathcal{L}(\mathbf{x}, \mathbf{v})$ in (5), the optimality condition for the primal update of PMM can be written as

$$\begin{aligned} \nabla f(\mathbf{x}_{t+1}) + (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{v}_t + \alpha(\mathbf{I} - \mathbf{Z})\mathbf{x}_{t+1} \\ + \epsilon(\mathbf{x}_{t+1} - \mathbf{x}_t) = \mathbf{0}. \end{aligned} \quad (51)$$

Further, notice that one of the KKT conditions of the optimization problem in (4) is

$$\nabla f(\mathbf{x}^*) + (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{v}^* = \mathbf{0}. \quad (52)$$

Moreover, the optimal solution $\mathbf{x}^* = [\tilde{\mathbf{x}}^*; \dots; \tilde{\mathbf{x}}^*]$ of (4) lies in $\text{null}\{\mathbf{I} - \mathbf{Z}\}$. Therefore, we obtain

$$\alpha(\mathbf{I} - \mathbf{Z})\mathbf{x}^* = \mathbf{0}. \quad (53)$$

Subtracting the equalities in (52) and (53) from (51) yields

$$\begin{aligned} \nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*) + (\mathbf{I} - \mathbf{Z})^{1/2}(\mathbf{v}_t - \mathbf{v}^*) \\ + \alpha(\mathbf{I} - \mathbf{Z})(\mathbf{x}_{t+1} - \mathbf{x}^*) + \epsilon(\mathbf{x}_{t+1} - \mathbf{x}_t) = \mathbf{0}. \end{aligned} \quad (54)$$

Regrouping the terms in (33) implies that \mathbf{v}_t is equivalent to

$$\mathbf{v}_t = \mathbf{v}_{t+1} - \alpha(\mathbf{I} - \mathbf{Z})^{1/2}(\mathbf{x}_{t+1} - \mathbf{x}^*). \quad (55)$$

Substituting \mathbf{v}_t in (54) by the expression in the right hand side of (55) leads to the claim in (34).

APPENDIX B PROOF OF THEOREM 1

According to Assumption 1, the global objective function f is strongly convex with constant m and its gradients ∇f are Lipschitz continuous with constant M . Considering

these assumptions, we obtain that the inner product $(\mathbf{x}_{t+1} - \mathbf{x}^*)^T (\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*))$ is lower bounded by

$$\begin{aligned} & \frac{mM}{m+M} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 + \frac{1}{m+M} \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 \\ & \leq (\mathbf{x}_{t+1} - \mathbf{x}^*)^T (\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)). \end{aligned} \quad (56)$$

The result in (34) shows that the difference $\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)$ is equal to $-(\mathbf{I} - \mathbf{Z})^{1/2}(\mathbf{v}_{t+1} - \mathbf{v}^*) - \epsilon(\mathbf{x}_{t+1} - \mathbf{x}_t)$. Apply this substitution into (56) and multiply both sides of the resulted inequality by 2 to obtain

$$\begin{aligned} & \frac{2mM}{m+M} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 + \frac{2}{m+M} \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 \\ & \leq -2(\mathbf{x}_{t+1} - \mathbf{x}^*)^T (\mathbf{I} - \mathbf{Z})^{1/2}(\mathbf{v}_{t+1} - \mathbf{v}^*) \\ & \quad - 2\epsilon(\mathbf{x}_{t+1} - \mathbf{x}^*)^T (\mathbf{x}_{t+1} - \mathbf{x}_t). \end{aligned} \quad (57)$$

Based on the result in (33), we can substitute $(\mathbf{x}_{t+1} - \mathbf{x}^*)^T (\mathbf{I} - \mathbf{Z})^{1/2}$ by $(1/\alpha)(\mathbf{v}_{t+1} - \mathbf{v}_t)^T$. Thus, we can rewrite (57) as

$$\begin{aligned} & \frac{2\alpha mM}{m+M} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 + \frac{2\alpha}{m+M} \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 \\ & \leq -2(\mathbf{v}_{t+1} - \mathbf{v}_t)^T (\mathbf{v}_{t+1} - \mathbf{v}^*) \\ & \quad - 2\alpha\epsilon(\mathbf{x}_{t+1} - \mathbf{x}^*)^T (\mathbf{x}_{t+1} - \mathbf{x}_t). \end{aligned} \quad (58)$$

Notice that for any vectors \mathbf{a} , \mathbf{b} , and \mathbf{c} we can write $2(\mathbf{a} - \mathbf{b})^T (\mathbf{a} - \mathbf{c}) = \|\mathbf{a} - \mathbf{b}\|^2 + \|\mathbf{a} - \mathbf{c}\|^2 - \|\mathbf{b} - \mathbf{c}\|^2$. By setting $\mathbf{a} = \mathbf{v}_{t+1}$, $\mathbf{b} = \mathbf{v}_t$, and $\mathbf{c} = \mathbf{v}^*$ we obtain that the inner product $2(\mathbf{v}_{t+1} - \mathbf{v}_t)^T (\mathbf{v}_{t+1} - \mathbf{v}^*)$ in (58) can be written as $\|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 + \|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2 - \|\mathbf{v}_t - \mathbf{v}^*\|^2$. Likewise, setting $\mathbf{a} = \mathbf{x}_{t+1}$, $\mathbf{b} = \mathbf{x}_t$, and $\mathbf{c} = \mathbf{x}^*$ implies that the inner product $2(\mathbf{x}_{t+1} - \mathbf{x}_t)^T (\mathbf{x}_{t+1} - \mathbf{x}^*)$ in (58) is equal to $\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 - \|\mathbf{x}_t - \mathbf{x}^*\|^2$. Hence, (58) can be simplified as

$$\begin{aligned} & \frac{2\alpha mM}{m+M} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 + \frac{2\alpha}{m+M} \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 \\ & \leq \alpha\epsilon \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \alpha\epsilon \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 - \alpha\epsilon \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \\ & \quad + \|\mathbf{v}_t - \mathbf{v}^*\|^2 - \|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2 - \|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2. \end{aligned} \quad (59)$$

Now using the definitions of the variable \mathbf{u} and matrix \mathcal{G} in (35) we can substitute $\|\mathbf{v}_t - \mathbf{v}^*\|^2 - \|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2 + \alpha\epsilon \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \alpha\epsilon \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2$ by $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2 - \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2$. Moreover, the squared norm $\|\mathbf{v}_{t+1} - \mathbf{v}_t\|^2$ is equivalent to $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|_{\alpha^2(\mathbf{I}-\mathbf{Z})}^2$ based on the result in (33). By applying these substitutions we can rewrite (59) as

$$\begin{aligned} & \frac{2\alpha mM}{m+M} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 + \frac{2\alpha}{m+M} \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 \\ & \leq \|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2 - \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2 - \alpha\epsilon \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ & \quad - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_{\alpha^2(\mathbf{I}-\mathbf{Z})}^2. \end{aligned} \quad (60)$$

Regrouping the terms in (60) leads to the following lower bound for the difference $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2 - \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2$,

$$\begin{aligned} & \|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2 - \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2 \\ & \geq \frac{2\alpha}{m+M} \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 + \alpha\epsilon \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ & \quad + \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_{\frac{2\alpha mM}{m+M} \mathbf{I} + \alpha^2(\mathbf{I}-\mathbf{Z})}^2. \end{aligned} \quad (61)$$

Observe that the result in (61) provides a lower bound for the decrement $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2 - \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2$. To prove the claim in (36), we need to show that for a positive constant δ we have $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2 - \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2 \geq \delta \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2$. Therefore, the inequality in (36) is satisfied if we can show that the lower bound in (61) is greater than $\delta \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2$ or equivalently

$$\begin{aligned} & \delta \|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2 + \delta\alpha\epsilon \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \\ & \leq \frac{2\alpha}{m+M} \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 + \alpha\epsilon \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ & \quad + \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_{\frac{2\alpha mM}{m+M} \mathbf{I} + \alpha^2(\mathbf{I}-\mathbf{Z})}^2. \end{aligned} \quad (62)$$

To prove that the inequality in (62) holds for some $\delta > 0$, we first find an upper bound for the squared norm $\|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2$ in terms of the summands in the right hand side of (62). To do so, consider the relation (34) along with the fact that \mathbf{v}_{t+1} and \mathbf{v}^* both lie in the column space of $(\mathbf{I} - \mathbf{Z})^{1/2}$. Note that there always exists a unique \mathbf{v}^* that lies in the column space of $(\mathbf{I} - \mathbf{Z})^{1/2}$ —check Lemma 1 in [28]. Since we know that both \mathbf{v}_{t+1} and \mathbf{v}^* lie in the column space of $(\mathbf{I} - \mathbf{Z})^{1/2}$, there exists a vector $\mathbf{r} \in \mathbb{R}^{np}$ such that $\mathbf{v}^* - \mathbf{v}_{t+1} = (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{r}$. This relation implies that $\|(\mathbf{I} - \mathbf{Z})^{1/2}(\mathbf{v}_{t+1} - \mathbf{v}^*)\|^2$ can be written as $\|(\mathbf{I} - \mathbf{Z})\mathbf{r}\|^2 = \mathbf{r}^T (\mathbf{I} - \mathbf{Z})^2 \mathbf{r}$. The eigenvalues of the matrix $(\mathbf{I} - \mathbf{Z})^2$ are the squared of eigenvalues of the matrix $(\mathbf{I} - \mathbf{Z})$. Thus, we can write $\mathbf{r}^T (\mathbf{I} - \mathbf{Z})^2 \mathbf{r} \geq \hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z}) \mathbf{r}^T (\mathbf{I} - \mathbf{Z}) \mathbf{r}$, where $\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})$ is the smallest non-zero eigenvalue of the matrix $\mathbf{I} - \mathbf{Z}$. Observing this inequality and the definition $\mathbf{v}^* - \mathbf{v}_{t+1} = (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{r}$ we can write

$$\|(\mathbf{I} - \mathbf{Z})^{1/2}(\mathbf{v}_{t+1} - \mathbf{v}^*)\|^2 \geq \hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z}) \|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2. \quad (63)$$

Moreover, from the inequality in (34) we obtain that $\|(\mathbf{I} - \mathbf{Z})^{1/2}(\mathbf{v}_{t+1} - \mathbf{v}^*)\|^2$ is bounded above by

$$\begin{aligned} & \|(\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{v}_{t+1} - \mathbf{v}^*\|^2 \\ & \leq \frac{\beta\epsilon^2}{(\beta-1)} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + \beta \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2, \end{aligned} \quad (64)$$

where $\beta > 1$ is a tunable free parameter. Replacing the norm $\|(\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{v}_{t+1} - \mathbf{v}^*\|^2$ in (64) by its lower bound in (63) follows that $\|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2$ is bounded above by

$$\begin{aligned} \|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2 & \leq \frac{\beta\epsilon^2}{(\beta-1)\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ & \quad + \frac{\beta}{\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})} \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2. \end{aligned} \quad (65)$$

Considering the result in (65) to satisfy the inequality in (62), which is a sufficient condition for the claim in (36), it remains

to show that

$$\begin{aligned}
& \frac{2\alpha}{m+M} \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 + \alpha\epsilon \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\
& + \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \frac{2\alpha m M}{m+M} \mathbf{I} + \alpha^2 (\mathbf{I} - \mathbf{Z}) \\
& \geq \frac{\delta\beta\epsilon^2}{(\beta-1)\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + \delta\epsilon\alpha \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \\
& + \frac{\delta\beta}{\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})} \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2. \tag{66}
\end{aligned}$$

To enable (66) and consequently enabling (62), we only need to verify that there exists $\delta > 0$ such that

$$\begin{aligned}
\frac{2\alpha m M}{m+M} \mathbf{I} + \alpha^2 (\mathbf{I} - \mathbf{Z}) & \succcurlyeq \delta\alpha\epsilon \mathbf{I}, \quad \frac{2\alpha}{m+M} \geq \frac{\delta\beta}{\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})}, \\
\alpha\epsilon & \geq \frac{\delta\beta\epsilon^2}{(\beta-1)\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})}. \tag{67}
\end{aligned}$$

The conditions in (67) are satisfied if the constant δ is chosen as in (37). Therefore, for δ in (37) the claim in (62) holds, which implies the claim in (36).

APPENDIX C PROOF OF LEMMA 2

Consider the primal update of ESOM in (14). By regrouping the terms we obtain that

$$\nabla f(\mathbf{x}_t) + (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{v}_t + \alpha(\mathbf{I} - \mathbf{Z})\mathbf{x}_t + \tilde{\mathbf{H}}_t(\mathbf{x}_{t+1} - \mathbf{x}_t) = \mathbf{0}, \tag{68}$$

where $\tilde{\mathbf{H}}_t$ is the inverse of the Hessian inverse approximation $\tilde{\mathbf{H}}_t^{-1}(K)$. Recall the definition of the exact Hessian \mathbf{H}_t in (9). Adding and subtracting the term $\mathbf{H}_t(\mathbf{x}_{t+1} - \mathbf{x}_t)$ to the expression in (68) yields

$$\begin{aligned}
& \nabla f(\mathbf{x}_t) + \nabla^2 f(\mathbf{x}_t)(\mathbf{x}_{t+1} - \mathbf{x}_t) + (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{v}_t \\
& + \alpha(\mathbf{I} - \tilde{\mathbf{Z}})\mathbf{x}_{t+1} + \epsilon(\mathbf{x}_{t+1} - \mathbf{x}_t) \\
& + (\tilde{\mathbf{H}}_t - \mathbf{H}_t)(\mathbf{x}_{t+1} - \mathbf{x}_t) = \mathbf{0}. \tag{69}
\end{aligned}$$

Now using the definition of the error vector \mathbf{e}_t in (40) we can rewrite (69) as

$$\begin{aligned}
& \nabla f(\mathbf{x}_{t+1}) + (\mathbf{I} - \mathbf{Z})^{1/2} \mathbf{v}_t \\
& + \alpha(\mathbf{I} - \tilde{\mathbf{Z}})\mathbf{x}_{t+1} + \epsilon(\mathbf{x}_{t+1} - \mathbf{x}_t) + \mathbf{e}_t = \mathbf{0}. \tag{70}
\end{aligned}$$

Notice that the result in (70) is identical to the expression for PMM in (51) except for the error term \mathbf{e}_t . To prove the claim in (39) from (70), it remains to follow the steps in (52)–(55).

APPENDIX D PROOF OF LEMMA 3

To prove the result in (42), we first use the result in Proposition 2 of [29]. It shows that when the eigenvalues of the Hessian $\nabla^2 f(\mathbf{x})$ are bounded above by M and the Hessian is Lipschitz

continuous with constant L we can write

$$\begin{aligned}
& \|\nabla f(\mathbf{x}_t) + \nabla^2 f(\mathbf{x}_t)(\mathbf{x}_{t+1} - \mathbf{x}_t) - \nabla f(\mathbf{x}_{t+1})\| \\
& \leq \|\mathbf{x}_{t+1} - \mathbf{x}_t\| \min \left\{ 2M, \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\| \right\}. \tag{71}
\end{aligned}$$

Considering the result in (71), it remains to find an upper bound for the second term of the error vector \mathbf{e}_t which is $(\tilde{\mathbf{H}}_t(K) - \mathbf{H}_t)(\mathbf{x}_{t+1} - \mathbf{x}_t)$. To do so, we develop first an upper bound for the norm $\|\tilde{\mathbf{H}}_t(K) - \mathbf{H}_t\|$. Notice that by factoring the term $\tilde{\mathbf{H}}_t(K)^{1/2}$ from left and right, and using the Cauchy-Schwarz inequality we obtain that

$$\|\tilde{\mathbf{H}}_t(K) - \mathbf{H}_t\| \leq \left\| \tilde{\mathbf{H}}_t(K)^{\frac{1}{2}} \right\|^2 \left\| \mathbf{I} - \tilde{\mathbf{H}}_t^{-\frac{1}{2}}(K) \mathbf{H}_t \tilde{\mathbf{H}}_t^{-\frac{1}{2}}(K) \right\|. \tag{72}$$

Note that the eigenvalues of the matrices $\mathbf{I} - \mathbf{H}_t \tilde{\mathbf{H}}_t^{-1}(K)$ and $\mathbf{I} - \tilde{\mathbf{H}}_t^{-1/2}(K) \mathbf{H}_t \tilde{\mathbf{H}}_t^{-1/2}(K)$ are the same since these two matrices are *similar*. In linear algebra, two matrices \mathbf{A} and $\tilde{\mathbf{A}}$ are called similar if $\tilde{\mathbf{A}} = \mathbf{P}^{-1} \mathbf{A} \mathbf{P}$ for an invertible matrix \mathbf{P} . Thus, we proceed to find bounds for the eigenvalues of $\mathbf{I} - \mathbf{H}_t \tilde{\mathbf{H}}_t^{-1}(K)$, to bound the norm in (72). According to Lemma 3 in [12], we can simplify $\mathbf{I} - \mathbf{H}_t \tilde{\mathbf{H}}_t^{-1}(K)$ as

$$\mathbf{I} - \mathbf{H}_t \tilde{\mathbf{H}}_t^{-1}(K) = (\mathbf{B} \mathbf{D}_t^{-1})^{K+1}. \tag{73}$$

Note that the matrices \mathbf{B} and \mathbf{D}_t in this paper are different from the ones in [12], but the analyses of them are very similar. Following the proof of Proposition 2 in [12], we define $\hat{\mathbf{D}} := 2\alpha(\mathbf{I} - \mathbf{Z}_d)$. Notice that the matrix $\hat{\mathbf{D}}$ is block diagonal where its i th diagonal block is $2\alpha(1 - w_{ii})\mathbf{I}_p$. Thus, $\hat{\mathbf{D}}$ is positive definite and invertible. Instead of studying an upper bound for the eigenvalues of $\mathbf{B} \mathbf{D}_t^{-1}$, we try to find an upper bound for the eigenvalues of its similar matrix $\mathbf{D}_t^{-1/2} \mathbf{B} \mathbf{D}_t^{-1/2}$ which is symmetric. We are allowed to write the product $\mathbf{D}_t^{-1/2} \mathbf{B} \mathbf{D}_t^{-1/2}$ as

$$\mathbf{D}_t^{-\frac{1}{2}} \mathbf{B} \mathbf{D}_t^{-\frac{1}{2}} = \left(\mathbf{D}_t^{-\frac{1}{2}} \hat{\mathbf{D}}^{\frac{1}{2}} \right) \left(\hat{\mathbf{D}}^{-\frac{1}{2}} \mathbf{B} \hat{\mathbf{D}}^{-1/2} \right) \left(\hat{\mathbf{D}}^{\frac{1}{2}} \mathbf{D}_t^{-\frac{1}{2}} \right). \tag{74}$$

The next step is to find an upper bound for the eigenvalues of $\mathbf{B} \hat{\mathbf{D}}^{-1}$ in (74). Based on the definitions of matrices \mathbf{B} and $\hat{\mathbf{D}}$, the product $\mathbf{B} \hat{\mathbf{D}}^{-1}$ is given by

$$\mathbf{B} \hat{\mathbf{D}}^{-1} = (\mathbf{I} - 2\mathbf{Z}_d + \mathbf{Z})(2(\mathbf{I} - \mathbf{Z}_d))^{-1}. \tag{75}$$

According to the result in Proposition 2 of [12], the eigenvalues of the matrix $(\mathbf{I} - 2\mathbf{Z}_d + \mathbf{Z})(2(\mathbf{I} - \mathbf{Z}_d))^{-1}$ are uniformly bounded by 0 and 1. Thus, we obtain that the eigenvalues of $\hat{\mathbf{D}}^{-1/2} \mathbf{B} \hat{\mathbf{D}}^{-1/2}$ are bounded by 0 and 1 and we can write

$$\|\hat{\mathbf{D}}^{-\frac{1}{2}} \mathbf{B} \hat{\mathbf{D}}^{-\frac{1}{2}}\| \leq 1. \tag{76}$$

According to the definitions of the matrices $\hat{\mathbf{D}}$ and \mathbf{D}_t , the product $\hat{\mathbf{D}}^{1/2} \mathbf{D}_t^{-1/2}$ is block diagonal and the i th diagonal block is given by

$$\left[\hat{\mathbf{D}} \mathbf{D}_t^{-1} \right]_{ii} = \left(\frac{\nabla^2 f_i(\mathbf{x}_{i,t}) + \epsilon \mathbf{I}}{2\alpha(1 - w_{ii})} + \mathbf{I} \right)^{-1}. \tag{77}$$

Based on Assumption 1, the eigenvalues of the local Hessians $\nabla^2 f_i(\mathbf{x}_i)$ are bounded by m and M . Further, notice that the diagonal elements w_{ii} of the weight matrix \mathbf{W} are bounded below

by c . Considering these bounds, we can show that the eigenvalues of the matrices $(1/2\alpha(1-w_{ii}))(\nabla^2 f_i(\mathbf{x}_{i,t}) + \epsilon \mathbf{I}) + \mathbf{I}$ for all $i = 1, \dots, n$ are bounded below by

$$\left[\frac{m + \epsilon}{2\alpha(1-c)} + 1 \right] \mathbf{I} \preceq \frac{\nabla^2 f_i(\mathbf{x}_{i,t}) + \epsilon \mathbf{I}}{2\alpha(1-w_{ii})} + \mathbf{I}. \quad (78)$$

By considering the bounds in (78), the eigenvalues of each block of the matrix $\hat{\mathbf{D}}\mathbf{D}_t^{-1}$, introduced in (77), are bounded above as

$$\left(\frac{\nabla^2 f_i(\mathbf{x}_{i,t}) + \epsilon \mathbf{I}}{2\alpha(1-w_{ii})} + \mathbf{I} \right)^{-1} \preceq \left[\frac{m + \epsilon}{2\alpha(1-c)} + 1 \right]^{-1} \mathbf{I}. \quad (79)$$

The upper bound in (79) for the eigenvalues of each diagonal block of the matrix $\hat{\mathbf{D}}\mathbf{D}_t^{-1}$ implies that the matrix norm $\|\hat{\mathbf{D}}\mathbf{D}_t^{-1}\|$ is bounded above by

$$\|\hat{\mathbf{D}}\mathbf{D}_t^{-1}\| \leq \rho := \frac{2\alpha(1-c)}{2\alpha(1-c) + m + \epsilon}. \quad (80)$$

Considering the upper bounds in (76) and (80) and the relation in (74) we obtain that

$$\|\mathbf{D}_t^{-\frac{1}{2}} \mathbf{B} \mathbf{D}_t^{-\frac{1}{2}}\| \leq \rho. \quad (81)$$

Thus, the eigenvalues of the positive definite symmetric matrix $\mathbf{D}_t^{-1/2} \mathbf{B} \mathbf{D}_t^{-1/2}$ are bounded by ρ . Hence, the eigenvalues of its similar matrix $\mathbf{B} \mathbf{D}_t^{-1}$ are bounded by ρ . This bound along with the result in (73) shows that the eigenvalues of the matrix $\mathbf{I} - \mathbf{H}_t \tilde{\mathbf{H}}_t^{-1}(K)$ are uniformly bounded by 0 and ρ^{K+1} . Therefore, the eigenvalues of its similar symmetric matrix $\mathbf{I} - \tilde{\mathbf{H}}_t^{-1/2}(K) \mathbf{H}_t \tilde{\mathbf{H}}_t^{-1/2}(K)$ are between 0 and ρ^K which implies that $\|\mathbf{I} - \tilde{\mathbf{H}}_t^{-1/2}(K) \mathbf{H}_t \tilde{\mathbf{H}}_t^{-1/2}(K)\| \leq \rho^{K+1}$. This result in conjunction with the inequality in (72) yields

$$\|\tilde{\mathbf{H}}_t(K) - \mathbf{H}_t\| \leq \rho^{K+1} \left\| \tilde{\mathbf{H}}_t(K)^{\frac{1}{2}} \right\|^2. \quad (82)$$

To bound the norm $\|\tilde{\mathbf{H}}_t(K)\|$, we first find a lower bound for the eigenvalues of the approximate Hessian inverse $\tilde{\mathbf{H}}_t^{-1}(K)$. Notice that according to the definition of the approximate Hessian inverse in (13), we can write

$$\tilde{\mathbf{H}}_t^{-1}(K) := \mathbf{D}_t^{-1} + \mathbf{D}_t^{-1} \sum_{u=1}^K (\mathbf{D}_t^{-1/2} \mathbf{B} \mathbf{D}_t^{-1/2})^u \mathbf{D}_t^{-1/2}. \quad (83)$$

Notice that according to the result in Proposition 1 of [12], the matrix $(\mathbf{I} - 2\mathbf{Z}_d + \mathbf{Z})$ is positive semidefinite which implies that $\mathbf{B} = \alpha(\mathbf{I} - 2\mathbf{Z}_d + \mathbf{Z})$ is also positive semidefinite. Thus, all the K summands in (83) are positive semidefinite and as a result we obtain that

$$\mathbf{D}_t^{-1} \preceq \tilde{\mathbf{H}}_t^{-1}(K). \quad (84)$$

The eigenvalues of $\mathbf{I} - \mathbf{Z}_d$ are bounded above by $1 - c$, since all the local weights w_{ii} are larger than c . This observation in conjunction with the strong convexity of the global objective function f implies that the eigenvalues of $\mathbf{D}_t = \nabla^2 f(\mathbf{x}_t) + \epsilon \mathbf{I} + 2\alpha(\mathbf{I} - \mathbf{Z}_d)$ are bounded above by $M + \epsilon + 2\alpha(1 - c)$. Therefore, the eigenvalues of \mathbf{D}_t^{-1} are bounded below as

$$\frac{1}{M + \epsilon + 2\alpha(1 - c)} \mathbf{I} \preceq \mathbf{D}_t^{-1}. \quad (85)$$

The results in (84) and (85) imply that the eigenvalues of the approximate Hessian inverse $\tilde{\mathbf{H}}_t^{-1}(K)$ are greater than $1/(M + \epsilon + 2\alpha(1 - c))$. Therefore, the eigenvalues of the positive definite matrix $\tilde{\mathbf{H}}_t(K)$ are smaller than $M + \epsilon + 2\alpha(1 - c)$ and we can write

$$\|\tilde{\mathbf{H}}_t(K)\| \leq M + \epsilon + 2\alpha(1 - c). \quad (86)$$

Considering the inequalities in (82) and (86) and using the Cauchy-Schwarz inequality we can show that the norm $\|(\tilde{\mathbf{H}}_t(K) - \mathbf{H}_t)(\mathbf{x}_{t+1} - \mathbf{x}_t)\|$ is bounded above by

$$\begin{aligned} & \left\| (\tilde{\mathbf{H}}_t(K) - \mathbf{H}_t)(\mathbf{x}_{t+1} - \mathbf{x}_t) \right\| \\ & \leq (M + \epsilon + 2\alpha(1 - c)) \rho^{K+1} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|. \end{aligned} \quad (87)$$

Observing the inequalities in (71) and (87) and using the triangle inequality the claim in (42) follows.

APPENDIX E PROOF OF THEOREM 2

Notice that in proving the claim in (44) we use some of the steps in the proof of Theorem 1 to avoid rewriting similar equations. First, note that according to the result in (39), the difference $\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)$ for the ESOM method can be written as

$$\begin{aligned} \nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*) &= -(\mathbf{I} - \mathbf{Z})^{1/2} (\mathbf{v}_{t+1} - \mathbf{v}^*) \\ &\quad - \epsilon(\mathbf{x}_{t+1} - \mathbf{x}_t) - \mathbf{e}_t. \end{aligned} \quad (88)$$

Now recall the the inequality in (56) and substitute the gradients difference $\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)$ in the inner product $(\mathbf{x}_{t+1} - \mathbf{x}^*)^T (\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*))$ by the expression in the right hand side of (88). Applying this substitution and multiplying both sides of the implied inequality by 2α follows

$$\begin{aligned} & \frac{2\alpha m M}{m + M} \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 + \frac{2\alpha}{m + M} \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 \\ & \leq -2\alpha \epsilon (\mathbf{x}_{t+1} - \mathbf{x}^*)^T (\mathbf{x}_{t+1} - \mathbf{x}_t) - 2\alpha (\mathbf{x}_{t+1} - \mathbf{x}^*)^T \mathbf{e}_t \\ & \quad - 2\alpha (\mathbf{x}_{t+1} - \mathbf{x}^*)^T (\mathbf{I} - \mathbf{Z})^{1/2} (\mathbf{v}_{t+1} - \mathbf{v}^*). \end{aligned} \quad (89)$$

By following the steps in (57)–(61), the result in (89) leads to a lower bound for $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2 - \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2$ as

$$\begin{aligned} & \|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2 - \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2 \\ & \geq \frac{2\alpha}{m + M} \|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 + \alpha \epsilon \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ & \quad + \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_{\frac{2\alpha m M}{m + M} \mathbf{I} + \alpha^2 (\mathbf{I} - \mathbf{Z})}^2 + 2\alpha (\mathbf{x}_{t+1} - \mathbf{x}^*)^T \mathbf{e}_t. \end{aligned} \quad (90)$$

Note that the inner product $2(\mathbf{x}_{t+1} - \mathbf{x}^*)^T \mathbf{e}_t$ is bounded below by $-(1/\zeta) \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 - \zeta \|\mathbf{e}_t\|^2$ for any positive constant $\zeta > 0$. Thus, the lower bound in (90) can be updated

as

$$\begin{aligned} & \|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2 - \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2 \\ & \geq \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_{\left(\frac{2\alpha m M}{m+M} - \frac{\alpha}{\zeta}\right)\mathbf{I} + \alpha^2(\mathbf{I} - \mathbf{Z})}^2 + \alpha\epsilon\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ & + \frac{2\alpha}{m+M}\|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 - \alpha\zeta\|\mathbf{e}_t\|^2. \end{aligned} \quad (91)$$

To establish (44), we need to show that the difference $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2 - \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2$ is bounded below by $\delta'_t\|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2$. To do so, we show that the lower bound for $\|\mathbf{u}_t - \mathbf{u}^*\|_{\mathcal{G}}^2 - \|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2$ in (91) is larger than $\delta'_t\|\mathbf{u}_{t+1} - \mathbf{u}^*\|_{\mathcal{G}}^2$, i.e.,

$$\begin{aligned} & \delta'_t\|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2 + \delta'_t\alpha\epsilon\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \\ & \leq \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_{\left(\frac{2\alpha m M}{m+M} - \frac{\alpha}{\zeta}\right)\mathbf{I} + \alpha^2(\mathbf{I} - \mathbf{Z})}^2 + \alpha\epsilon\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ & + \frac{2\alpha}{m+M}\|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 - \alpha\zeta\|\mathbf{e}_t\|^2. \end{aligned} \quad (92)$$

We proceed to find an upper bound for the squared norm $\|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2$ in terms of the summands in the right hand side of (92). Consider the relation (70) as well as the fact that \mathbf{v}_{t+1} and \mathbf{v}^* both lie in the column space of $(\mathbf{I} - \mathbf{Z})^{1/2}$. It follows that $\|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2$ is bounded above by

$$\begin{aligned} \|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2 & \leq \frac{\beta\epsilon^2}{(\beta-1)\hat{\lambda}}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + \frac{\beta\phi}{(\phi-1)\hat{\lambda}}\|\mathbf{e}_t\|^2 \\ & + \frac{\phi\beta}{\hat{\lambda}}\|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2, \end{aligned} \quad (93)$$

where we have used $\hat{\lambda}$ instead of $\hat{\lambda}_{\min}(\mathbf{I} - \mathbf{Z})$ to simplify notation. By substituting the upper bound in (93) for the squared norm $\|\mathbf{v}_{t+1} - \mathbf{v}^*\|^2$ in (92) we obtain a sufficient condition for the result in (92) which is given by

$$\begin{aligned} & \delta'_t\alpha\epsilon\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 + \frac{\delta'_t\beta\epsilon^2}{(\beta-1)\hat{\lambda}}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ & + \frac{\delta'_t\phi\beta}{\hat{\lambda}}\|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 + \frac{\delta'_t\beta\phi\alpha^2\|\mathbf{e}_t\|^2}{(\phi-1)\hat{\lambda}} \\ & \leq \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_{\left(\frac{2\alpha m M}{m+M} - \frac{\alpha}{\zeta}\right)\mathbf{I} + \alpha^2(\mathbf{I} - \mathbf{Z})}^2 + \alpha\epsilon\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ & + \frac{2\alpha}{m+M}\|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 - \alpha\zeta\|\mathbf{e}_t\|^2. \end{aligned} \quad (94)$$

Substitute the squared norm $\|\mathbf{e}_t\|^2$ terms in (94) by the upper bound in (42). It follows from this substitution and regrouping the terms that

$$\begin{aligned} 0 & \leq \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_{\left(\frac{2\alpha m M}{m+M} - \frac{\alpha}{\zeta} - \delta'_t\alpha\epsilon\right)\mathbf{I} + \alpha^2(\mathbf{I} - \mathbf{Z})}^2 \\ & + \left(\frac{2\alpha}{m+M} - \frac{\delta'_t\phi\beta}{\hat{\lambda}}\right)\|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2 \\ & + \left[\alpha\epsilon - \frac{\delta'_t\beta\epsilon^2}{(\beta-1)\hat{\lambda}} - \frac{\delta'_t\beta\phi\Gamma^2}{(\phi-1)\hat{\lambda}} - \alpha\zeta\Gamma^2\right]\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2. \end{aligned} \quad (95)$$

Notice that if the inequality in (95) is satisfied, then the result in (94) holds which implies the result in (92) and the linear convergence claim in (44). To satisfy the inequality in (95) we need

to make sure that the coefficients of the terms $\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2$, $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2$, and $\|\nabla f(\mathbf{x}_{t+1}) - \nabla f(\mathbf{x}^*)\|^2$ are non-negative. Therefore, the inequality in (95) holds if δ'_t satisfies

$$\begin{aligned} \frac{2\alpha m M}{m+M} - \frac{\alpha}{\zeta} - \delta'_t\alpha\epsilon & \geq 0, \quad \frac{2\alpha}{m+M} \geq \frac{\delta'_t\phi\beta}{\hat{\lambda}} \\ \alpha\epsilon & \geq \frac{\delta'_t\beta\epsilon^2}{(\beta-1)\hat{\lambda}} + \frac{\delta'_t\beta\phi\Gamma^2}{(\phi-1)\hat{\lambda}} + \alpha\zeta\Gamma^2. \end{aligned} \quad (96)$$

The conditions in (96) are satisfied if δ'_t is chosen as in (45). Thus, δ'_t in (45) satisfies the conditions in (96) and the claim in (44) holds.

APPENDIX F

The result in Theorem 2 holds if the interval $((m+M)/2mM, \epsilon/\Gamma_t^2)$ is non-empty or equivalently if the inequality $\epsilon > \Gamma_t^2(m+M)/2mM$ holds. However, Γ_t depends on ϵ which makes it unclear if there exists a choice of ϵ that satisfies the inequality $\epsilon > \Gamma_t^2(m+M)/2mM$. In the following proposition, we prove that the interval $((m+M)/2mM, \epsilon/\Gamma_t^2)$ is non-empty for a proper choice of ϵ .

Proposition 2: Consider ESOM as introduced in (8) (15). Recall the definition of Γ_t in (43). If the constant ϵ is chosen such that

$$\epsilon > \frac{m+M}{2mM} \left(2M + 2\alpha(1-c)\frac{M}{m}\right)^2, \quad (97)$$

then the inequality $\epsilon > \Gamma_t^2(m+M)/2mM$ holds and the set $((m+M)/2mM, \epsilon/\Gamma_t^2)$ is non-empty.

Proof: Note that the condition $\epsilon > \Gamma_t^2(m+M)/2mM$ is equivalent to

$$\Gamma_t < \frac{\sqrt{2\epsilon m M}}{\sqrt{m+M}}. \quad (98)$$

According to the definition of Γ_t , the expression $\rho := 2\alpha(1-c)/(2\alpha(1-c) + m + \epsilon)$, and the fact that $2M \geq \min\{2M, \frac{L}{2}\|\mathbf{x}_{t+1} - \mathbf{x}_t\|\}$, we can write

$$\Gamma_t \leq 2M + (M + \epsilon + 2\alpha(1-c)) \left(\frac{2\alpha(1-c)}{2\alpha(1-c) + m + \epsilon}\right)^{K+1}. \quad (99)$$

The results in (98) and (99) show that the inequality $\epsilon > \Gamma_t^2(m+M)/2mM$ holds if the following inequality holds,

$$\begin{aligned} 2M + (M + \epsilon + 2\alpha(1-c)) \left(\frac{2\alpha(1-c)}{2\alpha(1-c) + m + \epsilon}\right)^{K+1} \\ < \frac{\sqrt{2\epsilon m M}}{\sqrt{m+M}}. \end{aligned} \quad (100)$$

Thus, if the condition in (100) holds then we have $\epsilon > \Gamma_t^2(m+M)/2mM$. Note that $(2\alpha(1-c)/(2\alpha(1-c) + m + \epsilon))^{K+1} \leq 2\alpha(1-c)/(2\alpha(1-c) + m + \epsilon)$ for any $K \geq 0$. Thus, if the following inequality is satisfied the inequality in

(100) is also valid,

$$2M + (M + \epsilon + 2\alpha(1-c)) \left[\frac{2\alpha(1-c)}{2\alpha(1-c) + m + \epsilon} \right] < \frac{\sqrt{2\epsilon m M}}{\sqrt{m+M}}. \quad (101)$$

Considering that $m < M$ and $2\alpha(1-c) + \epsilon > 0$, we obtain that $(M + \epsilon + 2\alpha(1-c))/(m + \epsilon + 2\alpha(1-c)) \leq M/m$. Replacing $(M + \epsilon + 2\alpha(1-c))/(m + \epsilon + 2\alpha(1-c))$ in (101) by the upper bound M/m implies that

$$2M + 2\alpha(1-c) \frac{M}{m} < \frac{\sqrt{2\epsilon m M}}{\sqrt{m+M}}. \quad (102)$$

Note that if the condition in (102) holds then the condition in (101) is satisfied. The result in (102) shows that if ϵ satisfies

$$\epsilon > \frac{m+M}{2mM} \left(2M + 2\alpha(1-c) \frac{M}{m} \right)^2, \quad (103)$$

then the inequality in (102) and consequently the inequalities in (101) and (100) hold true which follows that the condition $\epsilon > \Gamma_t^2(m+M)/2mM$ is satisfied. ■

REFERENCES

- [1] F. Bullo, J. Cortés, and S. Martinez, *Distributed Control of Robotic Networks: A Mathematical Approach to Motion Coordination Algorithms*. Princeton, NJ, USA: Princeton Univ. Press, 2009.
- [2] Y. Cao, W. Yu, W. Ren, and G. Chen, "An overview of recent progress in the study of distributed multi-agent coordination," *IEEE Trans. Ind. Inf.*, vol. 9, pp. 427–438, Feb. 2013.
- [3] C. G. Lopes and A. H. Sayed, "Diffusion least-mean squares over adaptive networks: Formulation and performance analysis," *IEEE Trans. Signal Process.*, vol. 56, no. 7, pp. 3122–3136, Jul. 2008.
- [4] A. Ribeiro, "Ergodic stochastic optimization algorithms for wireless communication and networking," *IEEE Trans. Signal Process.*, vol. 58, no. 12, pp. 6369–6386, Dec. 2010.
- [5] A. Ribeiro, "Optimal resource allocation in wireless communication and networking," *EURASIP J. Wireless Commun. Netw.*, vol. 2012, no. 1, pp. 1–19, 2012.
- [6] I. D. Schizas, A. Ribeiro, and G. B. Giannakis, "Consensus in ad hoc WSNs with noisy links—part I: Distributed estimation of deterministic signals," *IEEE Trans. Signal Process.*, vol. 56, no. 1, pp. 350–364, Jan. 2008.
- [7] U. A. Khan, S. Kar, and J. M. Moura, "Diland: An algorithm for distributed sensor localization with noisy distance measurements," *IEEE Trans. Signal Process.*, vol. 58, no. 3, pp. 1940–1947, Mar. 2010.
- [8] M. Rabbat and R. Nowak, "Distributed optimization in sensor networks," in *Proc. ACM 3rd Int. Symp. Inf. Process. Sensor Netw.*, 2004, pp. 20–27.
- [9] R. Bekkerman, M. Bilenko, and J. Langford, *Scaling Up Machine Learning: Parallel and Distributed Approaches*. Cambridge, U.K.: Cambridge Univ. Press, 2011.
- [10] K. I. Tsianos, S. Lawlor, and M. G. Rabbat, "Consensus-based distributed optimization: Practical issues and applications in large-scale machine learning," in *Proc. 50th Annu. Allerton Conf. Commun. Control Comput.*, 2012, pp. 1543–1550.
- [11] V. Cevher, S. Becker, and M. Schmidt, "Convex optimization for big data: Scalable, randomized, and parallel algorithms for big data analytics," *IEEE Signal Process. Mag.*, vol. 31, no. 5, pp. 32–43, Sep. 2014.
- [12] A. Mokhtari, Q. Ling, and A. Ribeiro, "Network newton-part I: Algorithm and convergence," *arXiv:1504.06017*, 2015.
- [13] A. Mokhtari, Q. Ling, and A. Ribeiro, "Network newton-part II: Convergence rate and implementation," *arXiv:1504.06020*, 2015.
- [14] D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation: Numerical Methods*. Englewood Cliffs, NJ, USA: Prentice-Hall, 1989.
- [15] A. P. Ruszczyński, *Nonlinear Optimization*, vol. 13. Princeton, NJ, USA: Princeton Univ. Press, 2006.
- [16] M. G. Rabbat, R. D. Nowak, J. Bucklew, "Generalized consensus computation in networked systems with erasure links," in *Proc. IEEE 6th Workshop Signal Process. Adv. Wireless Commun.*, 2005, pp. 1088–1092.
- [17] M. R. Hestenes, "Multiplier and gradient methods," *J. Optim. Theory Appl.*, vol. 4, no. 5, pp. 303–320, 1969.
- [18] D. P. Bertsekas, *Constrained Optimization and Lagrange Multiplier Methods*. New York, NY, USA: Academic, 2014.
- [19] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Found. Trends Mach. Learn.*, vol. 3, no. 1, pp. 1–122, 2011.
- [20] W. Shi, Q. Ling, K. Yuan, G. Wu, and W. Yin, "On the linear convergence of the ADMM in decentralized consensus optimization," *IEEE Trans. Signal Process.*, vol. 62, no. 7, pp. 1750–1761, Jan. 2014.
- [21] N. Watanabe, Y. Nishimura, and M. Matsubara, "Decomposition in large system optimization using the method of multipliers," *J. Optim. Theory Appl.*, vol. 25, no. 2, pp. 181–193, 1978.
- [22] G. Stephanopoulos and A. W. Westerberg, "The use of Hestenes' method of multipliers to resolve dual gaps in engineering system optimization," *J. Optim. Theory Appl.*, vol. 15, no. 3, pp. 285–309, 1975.
- [23] J. M. Mulvey and A. Ruszczyński, "A diagonal quadratic approximation method for large scale linear programs," *Oper. Res. Lett.*, vol. 12, no. 4, pp. 205–215, 1992.
- [24] A. Ruszczyński, "On convergence of an augmented lagrangian decomposition method for sparse convex optimization," *Math. Oper. Res.*, vol. 20, no. 3, pp. 634–656, 1995.
- [25] R. Tappenden, P. Richtárik, and B. Büke, "Separable approximations and decomposition methods for the augmented Lagrangian," *Optim. Methods Softw.*, vol. 30, no. 3, pp. 643–668, 2015.
- [26] N. Chatzipanagiotis, D. Dentcheva, and M. M. Zavlanos, "An augmented lagrangian method for distributed optimization," *Math. Program.*, vol. 152, no. 1, pp. 405–434, Aug. 2015.
- [27] D. Jakovetic, J. Xavier, and J. M. Moura, "Cooperative convex optimization in networked systems: Augmented lagrangian algorithms with directed gossip communication," *IEEE Trans. Signal Process.*, vol. 59, no. 8, pp. 3889–3902, Feb. 2011.
- [28] Q. Ling, W. Shi, G. Wu, and A. Ribeiro, "DLM: Decentralized linearized alternating direction method of multipliers," *IEEE Trans. Signal Process.*, vol. 63, no. 15, pp. 4051–4064, Aug. 2015.
- [29] A. Mokhtari, W. Shi, Q. Ling, and A. Ribeiro, "DQM: Decentralized quadratically approximated alternating direction method of multipliers," *IEEE Trans. Signal Process.*, vol. 64, no. 19, pp. 5158–5173, Oct. 2016.
- [30] M. Zargham, A. Ribeiro, A. Ozdaglar, and A. Jadbabaie, "Accelerated dual descent for network flow optimization," *IEEE Trans. Autom. Control*, vol. 59, no. 4, pp. 905–920, Apr. 2014.
- [31] A. Mokhtari, Q. Ling, and A. Ribeiro, "An approximate newton method for distributed optimization," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2015, pp. 2959–2963.
- [32] W. Shi, Q. Ling, G. Wu, and W. Yin, "Extra: An exact first-order algorithm for decentralized consensus optimization," *SIAM J. Optim.*, vol. 25, no. 2, pp. 944–966, 2015.
- [33] S. Boyd, P. Diaconis, and L. Xiao, "Fastest mixing markov chain on a graph," *SIAM Rev.*, vol. 46, no. 4, pp. 667–689, 2004.
- [34] D. Bajovic, D. Jakovetic, N. Krejic, and N. K. Jerinkic, "Newton-like method with diagonal correction for distributed optimization," *arXiv:1509.01703*, 2015.
- [35] A. Nedic and A. Ozdaglar, "Distributed subgradient methods for multi-agent optimization," *IEEE Trans. Autom. Control*, vol. 54, no. 1, pp. 48–61, Jan. 2009.



Aryan Mokhtari received the B.Sc. degree in electrical engineering from Sharif University of Technology, Tehran, Iran, in 2011, and the M.S. degree in electrical engineering from the University of Pennsylvania, Philadelphia, PA, USA, in 2014. Since 2012, he has been working toward the Ph.D. degree in the Department of Electrical and Systems Engineering, University of Pennsylvania, Philadelphia, PA. From June to August 2010, he was an intern at the Advanced Digital Sciences Center, Singapore. He was a research intern with the Big-Data Machine Learning Group, Yahoo!, Sunnyvale, CA, from June to August 2016. His research interests include the areas of optimization, machine learning, control, and signal processing. His current research focuses on developing stochastic, distributed (parallel), and decentralized methods for large-scale optimization problems.



Wei Shi received the B.E. degree in automation and the Ph.D. degree in control science and engineering both from the University of Science and Technology of China, Hefei, China, in 2010 and 2015, respectively. From 2015 to 2016 he was a Postdoctoral Research Associate at Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, Urbana. He is currently a Postdoctoral Research Associate at Boston University, Boston, MA, USA. His research interests include optimization and its applications in signal processing and control.



Qing Ling received the B.E. degree in automation and the Ph.D. degree in control theory and control engineering both from the University of Science and Technology of China, Hefei, China, in 2001 and 2006, respectively. From 2006 to 2009, he was a Postdoctoral Research Fellow in the Department of Electrical and Computer Engineering, Michigan Technological University. Since 2009, he has been an Associate Professor in the Department of Automation, University of Science and Technology of China. His current research interests include the decentralized optimization

of networked multiagent systems.



Alejandro Ribeiro received the B.Sc. degree in electrical engineering from the Universidad de la Republica Oriental del Uruguay, Montevideo, Uruguay, in 1998 and the M.Sc. and Ph.D. degree in electrical engineering from the Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN, USA, in 2005 and 2007, respectively. From 1998 to 2003, he was a Member of the Technical Staff at Bellsouth Montevideo. After the M.Sc. and Ph.D. degrees, in 2008

he joined the University of Pennsylvania, Philadelphia, PA, USA, where he is currently the Rosenbluth Associate Professor at the Department of Electrical and Systems Engineering. His research interests include the applications of statistical signal processing to the study of networks and networked phenomena. His focus is on structured representations of networked data structures, graph signal processing, network optimization, robot teams, and networked control. Dr. Ribeiro received the 2014 O. Hugo Schuck Best Paper Award, the 2012 S. Reid Warren, Jr. Award presented by Penn's undergraduate student body for outstanding teaching, the NSF CAREER Award in 2010, and Paper Awards at the 2016 SSP Workshop, 2016 SAM Workshop, 2015 Asilomar SSC Conference, ACC 2013, ICASSP 2006, and ICASSP 2005. He is a Fulbright Scholar and a Penn Fellow.