

1 **NONLINEAR CONSENSUS+INNOVATIONS UNDER CORRELATED**
2 **HEAVY-TAILED NOISES: MEAN SQUARE CONVERGENCE RATE**
3 **AND ASYMPTOTICS**

4 MANOJLO VUKOVIC*, DUSAN JAKOVETIC†, DRAGANA BAJOVIC‡, AND SOUMMYA
5 KAR§

6 **Abstract.** We consider distributed recursive estimation of consensus+innovations type in the
7 presence of heavy-tailed sensing and communication noises. We allow that the sensing and commu-
8 nication noises are mutually correlated while independent identically distributed (i.i.d.) in time, and
9 that they may both have infinite moments of order higher than one (hence having infinite variances).
10 Such heavy-tailed, infinite-variance noises are highly relevant in practice and are shown to occur,
11 e.g., in dense internet of things (IoT) deployments. We develop a consensus+innovations distributed
12 estimator that employs a general nonlinearity in both consensus and innovations steps to combat
13 the noise. We establish the estimator’s almost sure convergence, asymptotic normality, and mean
14 squared error (MSE) convergence. Moreover, we establish and explicitly quantify for the estimator
15 a sublinear MSE convergence rate. We then quantify through analytical examples the effects of
16 the nonlinearity choices and the noises correlation on the system performance. Finally, numerical
17 examples corroborate our findings and verify that the proposed method works in the simultaneous
18 heavy-tail communication-sensing noise setting, while existing methods fail under the same noise
19 conditions.

20 **Key words.** nonlinear mappings, consensus+innovations, distributed estimation, heavy-tailed
21 noise, mean square convergence rate, correlated noises

22 **AMS subject classifications.** 93E10, 93E35, 60G35, 94A13, 62M05

23 **1. Introduction.** We consider a distributed estimation problem where a net-
24 work of agents cooperates to estimate an unknown static vector parameter $\theta^* \in$
25 \mathbb{R}^M . Specifically, we are interested in *consensus+innovations* distributed estimation,
26 e.g., [18, 16, 17]. With consensus+innovations, each agent iteratively updates its un-
27 known parameter’s estimate by 1) exchanging its estimate with immediate neighbors
28 in the network; and 2) assimilating a newly acquired observation (measurement).

29 Consensus+innovations distributed estimators have been extensively studied, e.g.,
30 [18, 16, 17]; see also [20, 22, 23, 27, 30, 24, 38] for related diffusion-type and other
31 methods. Typically, such distributed estimators exhibit strong convergence guaran-
32 tees under various imperfection models (noises) in 1) sensing (observations) and/or
33 2) inter-agent communications. For example, reference [18] establishes almost sure
34 (a.s.) convergence and asymptotic normality of the estimators developed therein.
35 The authors of [18] allow for an observation noise with finite variance and a network
36 model that accounts for random link failures and dithered quantization (effectively an
37 additive noise with finite variance). Reference [16] considers consensus+innovations
38 distributed estimation in the presence of random link failures without quantization or
39 additive noise, and it develops estimators that are asymptotically efficient, i.e., that
40 achieve the minimal possible asymptotic variance. The authors of [17] propose adap-

*University of Novi Sad, Faculty of Technical Sciences, Department of Fundamental Sciences (manojlo.vukovic@uns.ac.rs).

†University of Novi Sad, Faculty of Sciences, Department of Mathematics and Informatics (dusan.jakovetic@dmi.uns.ac.rs).

‡University of Novi Sad, Faculty of Technical Sciences, Department of Power, Electronic and Communication Engineering (dbajovic@uns.ac.rs).

§Department of Electrical and Computer Engineering, Carnegie Mellon University (soumyak@andrew.cmu.edu).

41 tive asymptotically efficient estimators, wherein the innovation gains are adaptively
 42 learned during the algorithm progress. Consensus+innovations distributed detection
 43 and related distributed detection methods have also been considered, e.g., [25, 3, 2, 14].
 44 The above distributed estimation and distributed detection-related works typically as-
 45 sume that the noises have finite moments of a certain order greater than two, and
 46 hence they have finite variance.

47 It is highly relevant to investigate distributed estimators in the presence of heavy-
 48 tailed *communication* and *sensing* noises, as they arise in many application scenarios.
 49 For example, edge devices in Internet of Things (IoT) systems or sensor networks can
 50 be subject to noise distributions that may not have finite moments of order higher
 51 than one, e.g., [6, 31, 12, 37, 11, 7], like, e.g., symmetric α -stable noise distributions.
 52 This effect may occur due to interference, e.g., when wireless sensor network is rela-
 53 tively densely deployed. In this case, the signals of neighboring nodes interfere with
 54 each other and corrupt the signal to be received. References [10, 36] analyze the prob-
 55 ability distribution of the interference and demonstrate that it has heavy-tails. More
 56 precisely, [10, 36] show that the interference power has an alpha-stable distribution
 57 in a network with infinite radius and no guard zone when the interferers are placed
 58 according to a Poisson point process, where alpha depends on the path loss coefficient
 59 between the interferers and the receiver (see [10, 36] for details). Empirical evidence
 60 for the emergence of heavy-tail interference noise in certain IoT systems has been
 61 provided in [6].

62 Moreover, observation and communication noises may be mutually correlated
 63 due to the common interference processes in the environment that the sensing and
 64 communication devices are exposed to.

65 Several recent works [19, 21, 35, 33, 5, 1, 4, 26] consider distributed estimation
 66 methods in the presence of *impulsive observations noise*,¹ but still assuming a *finite*
 67 *noise variance* and *no communication noise*. For example, reference [19] introduces a
 68 method based on Wilcoxon-norm; [21] utilizes a Huber-loss function; and [35] adopts a
 69 mean error minimization approach. Robust distributed estimation methods based on
 70 adaptive subgradient projections are considered in [33, 5]. To cope with the impulsive
 71 observation noise, several references employ a certain *nonlinearity* in the innovation
 72 step. Reference [1] develops a method that adaptively learns an optimized nonlinearity
 73 at the innovation step for each agent in the network. Reference [4] employs a satura-
 74 tion nonlinearity in the innovation step to cope with measurement attacks. Further
 75 results on distributed estimation under impulsive observations noise can be found in a
 76 recent survey [26]. Very recently, we have developed a consensus+innovations distrib-
 77 uted estimator [15] that provably works under a heavy-tailed communications noise
 78 and a light-tailed observations noise. Specifically, under the assumed setting, [15]
 79 establishes almost sure convergence and asymptotic normality of the method therein.
 80 However, [15] is not concerned with mean squared error (MSE) rate analysis of the
 81 method. While asymptotic normality is a useful result that provides the algorithm's
 82 rate of convergence (in the weak convergence sense) *asymptotically*, it does not capture
 83 the (MSE) algorithm behavior in non-asymptotic regimes.

84 In summary, we identify for the current literature the following major gaps with
 85 respect to design and analysis of distributed estimation methods under heavy-tailed

¹As explained in, e.g., [1], an impulsive noise may be described as one whose realizations contain sparse, random samples of amplitude much higher than nominally accounted for. Impulsive noise may have a finite or infinite variance. Existing works on distributed estimation in impulsive noises assume a *finite noise variance*.

noises. 1) All existing works assume a finite observations noise variance. That is, even when impulsive observation noise is assumed, existing works still require the variance of the noise to be finite. This assumption can be restrictive and is violated for several commonly used heavy-tail noise models like α -stable distributions [11]. 2) No existing work simultaneously handles heavy-tailed (infinite-variance) sensing and heavy-tailed (infinite-variance) observation noises. 3) MSE convergence rate analysis has not been developed for distributed estimation in the presence of either infinite-variance sensing and/or infinite-variance communication noises. 4) Existing works on distributed estimation in the presence of infinite-variance (either sensing and/or communication noises) assume mutually independent sensing and communication noises.

Contributions. In this paper, we close the gaps identified above by developing a nonlinear consensus+innovations distributed estimator that provably works under the simultaneous presence of correlated heavy-tailed (infinite variance) observation and communication noises. We allow for a very general model of the sensing and communication noises, only assuming that they exhibit symmetric zero-mean distributions with finite first moments. Hence, the variances of both sensing and communication noises may be infinite. Moreover, we allow that, for a fixed time instant t , the additive sensing and communication noises may be mutually dependent, while they are both independent identically distributed (i.i.d.) in time. The proposed estimator employs a generic nonlinearity both at the innovations and the consensus terms. The encompassed nonlinearities are very general and include a broad class of (possibly discontinuous) odd functions, such as the component-wise sign and clipping functions. We establish for the proposed estimator almost sure convergence, asymptotic normality, and we explicitly evaluate the corresponding asymptotic variance. Furthermore, we establish for the proposed method, under a carefully designed step size sequence, a MSE convergence rate $O(1/t^\kappa)$, and we quantify the rate $\kappa \in (0, 1)$ in terms of the system parameters. In addition, we quantify through analytical examples the effects of correlation between sensing and observation noises, and we demonstrate how the derived asymptotic covariance results may be used as a guideline to optimize the employed nonlinearities for a problem at hand. Finally, we compare the proposed method with existing works in [1] and [15], both through analytical examples and by simulation. Most notably, we show that the existing methods fail to converge under the simultaneous presence of heavy-tailed (infinite-variance) observation and communication noises, while the proposed method provably works in the heavy-tailed setting.

Paper organization. Section 2 provides a description of the distributed estimation model that is considered and also gives all basic assumptions. In Section 3, we present the proposed nonlinear consensus+innovations estimator. Section 4 establishes almost sure convergence, asymptotic normality and the MSE rate of the proposed distributed estimator. Section 5 presents analytical and numerical examples. The conclusion is given in Section 6. Some auxiliary supporting arguments are provided in [34].

Notation. We denote by \mathbb{R} the set of real numbers and by \mathbb{R}^m the m -dimensional Euclidean real coordinate space. We use normal lower-case letters for scalars, lower case boldface letters for vectors, and upper case boldface letters for matrices. Further, to represent a vector $\mathbf{a} \in \mathbb{R}^m$ through its component, we write $\mathbf{a} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m]^\top$ and we denote by: \mathbf{a}_i or $[\mathbf{a}_i]$, as appropriate, the i -th element of vector \mathbf{a} ; \mathbf{A}_{ij} or $[\mathbf{A}_{ij}]$, as appropriate, the entry in the i -th row and j -th column of a matrix \mathbf{A} ; \mathbf{A}^\top the transpose of a matrix \mathbf{A} ; \otimes the Kronecker product of matrices. Further, we use either $\mathbf{a}^\top \mathbf{b}$ or $\langle \mathbf{a}, \mathbf{b} \rangle$ for the inner products of vectors \mathbf{a} and \mathbf{b} . Next, we let \mathbf{I} , $\mathbf{0}$, and

136 $\mathbf{1}$ be, respectively, the identity matrix, the zero vector, and the column vector with
 137 unit entries; $\text{Diag}(\mathbf{a})$ the diagonal matrix whose diagonal entries are the elements of
 138 vector \mathbf{a} ; \mathbf{J} the $N \times N$ matrix $\mathbf{J} := (1/N)\mathbf{1}\mathbf{1}^\top$. When appropriate, we indicate the
 139 matrix or vector dimension through a subscript. Next, $\mathbf{A} \succ 0$ ($\mathbf{A} \succeq 0$) means that
 140 the symmetric matrix A is positive definite (respectively, positive semi-definite). We
 141 further denote by: $\|\cdot\| = \|\cdot\|_2$ the Euclidean (respectively, spectral) norm of its vector
 142 (respectively, matrix) argument; $\lambda_i(\cdot)$ the i -th smallest eigenvalue; $g'(v)$ the derivative
 143 evaluated at v of a function $g : \mathbb{R} \rightarrow \mathbb{R}$; $\nabla h(\mathbf{w})$ and $\nabla^2 h(\mathbf{w})$ the gradient and Hessian,
 144 respectively, evaluated at w of a function $h : \mathbb{R}^m \rightarrow \mathbb{R}$, $m > 1$; $\mathbb{P}(\mathcal{A})$ and $\mathbb{E}[u]$ the
 145 probability of an event \mathcal{A} and expectation of a random variable u , respectively; and
 146 by $\text{sign}(a)$ the sign function, i.e., $\text{sign}(a) = 1$, for $a > 0$, $\text{sign}(a) = -1$, for $a < 0$, and
 147 $\text{sign}(0) = 0$. Finally, for two positive sequences η_n and χ_n , we have: $\eta_n = O(\chi_n)$ if
 148 $\limsup_{n \rightarrow \infty} \frac{\eta_n}{\chi_n} < \infty$.

149 **2. Problem model and basic assumptions.** We consider a network of N
 150 agents (sensors), through which the parameter of interest $\boldsymbol{\theta}^* \in \mathbb{R}^M$ is to be estimated.
 151 At each time $t = 0, 1, \dots$, each agent $i = 1, 2, \dots, N$ observes parameter $\boldsymbol{\theta}^*$ following
 152 the linear regression model:

$$153 \quad (2.1) \quad z_i^t = \mathbf{h}_i^\top \boldsymbol{\theta}^* + n_i^t.$$

154 Here, $z_i^t \in \mathbb{R}$ is the observation, $\mathbf{h}_i \in \mathbb{R}^M$ is the deterministic, non-zero regression
 155 vector known only by agent i and $n_i^t \in \mathbb{R}$ is the observation noise. The underlying
 156 topology is modeled via a graph $G = (V, E)$, where $V = \{1, \dots, N\}$ is the set of agents
 157 and E is the set of links, i.e., $\{i, j\} \in E$ if there exists a link between agents i and j .
 158 We also define the set of all arcs E_d in the following way: if $\{i, j\} \in E$ then $(i, j) \in E_d$
 159 and $(j, i) \in E_d$. We denote by $\Omega_i = \{j \in V : \{i, j\} \in E\}$ set of neighbors of agent i
 160 (excluding i) and by $\mathbf{D} = \text{Diag}(\{d_i\})$ the degree matrix, where $d_i = |\Omega_i|$ is the number
 161 of neighbors of agent i . The graph Laplacian matrix \mathbf{L} is defined by $\mathbf{L} = \mathbf{D} - \mathbf{A}$,
 162 where \mathbf{A} is the adjacency matrix, which is a zero-one symmetric matrix with zero
 163 diagonal, such that, for $i \neq j$, $\mathbf{A}_{ij} = 1$ if and only if $\{i, j\} \in E$. Let us denote by
 164 $(\Omega, \mathcal{F}, \mathbb{P})$ the underlying probability space.
 165 We make the following assumptions.

167 **Assumption 2.1. Network model and Observability:**

- 168 1. Graph $G = (V, E)$ is undirected, simple (no self or multiple links) and static;
- 169 2. The matrix $\sum_{i=1}^N \mathbf{h}_i \mathbf{h}_i^\top$ is invertible;

170 The condition 2 in Assumption 2.1 ensures that (2.1) is observable, i.e., a centralized
 171 estimator (e.g., least squares) that collects all $z_i^t, i = 1, 2, \dots, N$, for all t , and has
 172 knowledge of all vectors $\mathbf{h}_i, i = 1, 2, \dots, N$, is consistent.

173 **Assumption 2.2. Observation noise:**

- 174 1. For each agent $i = 1, \dots, N$, the observation noise sequence $\{n_i^t\}$ in (2.1), is
 175 independent identically distributed (i.i.d.);
- 176 2. At each agent $i = 1, \dots, N$ at each time $t = 0, 1, \dots$, noise n_i^t has the same
 177 probability density function p_o .
- 178 3. Random variables n_i^t and n_j^s are mutually independent whenever the tuple
 179 (i, t) is different from (j, s) ;
- 180 4. The pdf p_o is symmetric, i.e. $p_o(u) = p_o(-u)$, for every $u \in \mathbb{R}$, and $p_o(u) > 0$
 181 for $|u| \leq c_o$, for some constant $c_o > 0$;
- 182 5. There holds that with $\int |u| p_o(u) du < \infty$.

183 If there is an arc between agents i and j , i.e., $(i, j) \in E_d$, we denote by $\boldsymbol{\xi}_{ij}^t$ communi-
 184 cation noise that is injected when agent j communicates to agent i at time instant t

185 (see ahead algorithm (3.1)).

186 *Assumption 2.3. Communication noise:*

- 187 1. Additive communication noise $\{\xi_{ij}^t\}, \xi_{ij}^t \in \mathbb{R}^M$ is i.i.d. in time t , and inde-
- 188 pendent across different arcs $(i, j) \in E_d$.
- 189 2. Each random variable $[\xi_{ij}^t]_\ell$, for each $t = 0, 1, \dots$, for each arc (i, j) , for each
- 190 entry $\ell = 1, \dots, M$, has the same probability density function p_c .
- 191 3. The pdf p_c is symmetric, i.e. $p_c(u) = p_c(-u)$, for every $u \in \mathbb{R}$ and $p_c(u) > 0$
- 192 for $|u| \leq c_c$, for some constant $c_c > 0$;
- 193 4. There holds that $\int |u|p_c(u)du < \infty$.

194 *Remark 2.4.* Notice here that from the symmetry of the probability density func-

195 tions p_o and p_c , it follows that both of the distributions are zero mean. Moreover,

196 notice that we do not assume that observation and communication noises are mutually

197 independent for a fixed t . However, they are both i.i.d. in time.

198 *Remark 2.5.* Condition 2 in Assumptions 2.2 and 2.3 can be relaxed in the sense

199 that it can be assumed that \mathbf{n}^t has joint probability density function p_o and ξ_{ij}^t has

200 the joint probability density function $p_{c,ij}$. (see Appendix C in [34]). The reason why

201 there is condition 4 in the Assumption 2.2 and condition 3 in the Assumption 2.3 will

202 become clear later.

203 For future reference, a compact vector form of (2.1) is:

204 (2.2)
$$\mathbf{z}^t = \mathbf{H}(\mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \mathbf{n}^t,$$

206 where, $\mathbf{z}^t = [z_1^t, z_2^t, \dots, z_N^t]^\top \in \mathbb{R}^N$ is the observation vector, $\mathbf{H} \in \mathbb{R}^{N \times (MN)}$ is the

207 regression matrix whose i -th row vector equals $[\mathbf{0}, \dots, \mathbf{0}, \mathbf{h}_i^\top, \mathbf{0}, \dots, \mathbf{0}] \in \mathbb{R}^{MN}$, where the

208 i -th block of size M equals \mathbf{h}_i^\top , and the other M -size blocks are the zero vectors; and

209 $\mathbf{n}^t = [n_1^t, n_2^t, \dots, n_N^t]^\top \in \mathbb{R}^N$ is the noise vector at time t .

210 **3. Proposed algorithm.** In order to estimate the unknown parameter $\boldsymbol{\theta}^* \in$

211 \mathbb{R}^M , in the presence of heavy-tailed observation noise and heavy-tailed communica-

212 tion noise, each agent uses a nonlinear consensus+innovations strategy. Therein, the

213 impact of the two heavy-tailed noises is mitigated by nonlinearities that have been

214 added to both consensus and innovation steps.

215 In more detail, each agent i at each time $t = 0, 1, \dots$, generates a sequence of estimates

216 $\{\mathbf{x}_i^t\}_{t \geq 0}$ of unknown parameter $\boldsymbol{\theta}^*$ by the following algorithm:

217 (3.1)
$$\mathbf{x}_i^{t+1} = \mathbf{x}_i^t - \alpha_t \left(\frac{b}{a} \sum_{j \in \Omega_i} \Psi_c(\mathbf{x}_i^t - \mathbf{x}_j^t + \xi_{ij}^t) - \mathbf{h}_i \Psi_o(z_i^t - \mathbf{h}_i^\top \mathbf{x}_i^t) \right).$$

218 Here, α_t is a step-size, and $a, b > 0$ are constants. We consider a family of decaying

219 step-size choices $\alpha_t = a/(t+1)^\delta$, $\delta \in (0.5, 1]$. As shown later, the step-size (values of

220 a and δ) should be designed appropriately in order for good properties (e.g., a.s. con-

221 vergence, MSE rate guarantees) of the algorithm to hold. Functions $\Psi_o : \mathbb{R} \rightarrow \mathbb{R}$ and

222 $\Psi_c : \mathbb{R}^M \rightarrow \mathbb{R}^M$ are non-linear functions and function Ψ_c operates component-wise by

224 abusing notation, i.e., for $\mathbf{y} \in \mathbb{R}^M$, we set that $\Psi_c(\mathbf{y}) = [\Psi_c(y_1), \Psi_c(y_2), \dots, \Psi_c(y_M)]$.

225 Also, functions Ψ_c and Ψ_o satisfy Assumption 3.1. We compare the proposed method

226 (3.1) with the \mathcal{LU} scheme in [18] and the scheme in [15]. Compared with these

227 schemes, (3.1) introduces a nonlinearity in the innovation step as well. \mathcal{LU} is obtained

228 from (3.1) by setting both of the nonlinearities Ψ_c and Ψ_o to identity functions and

229 $\delta = 1$, the method in [15] is recovered from (3.1) by setting Ψ_o to the identity function

230 and $\delta = 1$.

231 **Assumption 3.1. Nonlinearity Ψ :**

232 The non-linear function $\Psi : \mathbb{R} \rightarrow \mathbb{R}$ satisfies the following properties:

- 233 1. Function Ψ is odd, i.e., $\Psi(a) = -\Psi(-a)$, for any $a \in \mathbb{R}$;
 234 2. $\Psi(a) > 0$, for any $a > 0$.
 235 3. Function Ψ is a monotonically nondecreasing function;
 236 4. Ψ is continuous, except possibly on a point set with Lebesgue measure of
 237 zero. Moreover, Ψ is piecewise differentiable;
 238 5. $|\Psi(a)| \leq c_1$, for some constant $c_1 > 0$.
 239 6. Ψ is either discontinuous at zero, or $\Psi(u)$ is strictly increasing for $u \in$
 240 $(-c_2, c_2)$, for some $c_2 > 0$.

241 As it will become clear ahead, the role of Ψ_c and Ψ_o is to lower the impact of the
 242 heavy-tailed noise that occurs in the regression model and in the communication
 243 between agents. As it is presented in [15], there are many nonlinear functions which
 244 satisfy Assumption 3.1. Now, we add more assumptions on the observation and
 245 communication noises through the following assumption.

246 At each time $t = 0, 1, \dots$, a compact vector form of algorithm (3.1) is

247 (3.2)
$$\mathbf{x}^{t+1} = \mathbf{x}^t - \alpha_t \left(\frac{b}{a} \mathbf{L}_{\Psi_c}(\mathbf{x}) - \mathbf{H}^\top \Psi_o(\mathbf{z}^t - \mathbf{H}\mathbf{x}^t) \right).$$

248 Here, $\mathbf{x}^t = [\mathbf{x}_1^t, \mathbf{x}_2^t, \dots, \mathbf{x}_N^t]^\top \in \mathbb{R}^{MN}$, map $\mathbf{L}_{\Psi_c}(\mathbf{x}) : \mathbb{R}^{MN} \rightarrow \mathbb{R}^{MN}$ is defined by

249
$$\mathbf{L}_{\Psi_c}(\mathbf{x}) = \begin{bmatrix} \vdots \\ \sum_{j \in \Omega_i} \Psi_c(\mathbf{x}_i - \mathbf{x}_j + \boldsymbol{\xi}_{ij}) \\ \vdots \end{bmatrix},$$

250 where, the blocks $\sum_{j \in \Omega_i} \Psi_c(\mathbf{x}_i - \mathbf{x}_j + \boldsymbol{\xi}_{ij}) \in \mathbb{R}^M$ are stacked one on top of another for
 251 $i = 1, \dots, N$.

252 **4. Theoretical results.** In subsection 4.1 we express algorithm (3.2) in more
 253 general way, that will be used in the following subsections. Subsection 4.2 presents
 254 the statement and the proof of almost sure convergence of algorithm (3.1). In sub-
 255 section 4.3 we state and prove asymptotic normality and calculate the corresponding
 256 asymptotic variance. Subsection 4.4 presents and proves results on MSE rates.

257 **4.1. Setting up analysis.** In this subsection we rewrite algorithm (3.1) in the
 258 form suitable for stating the main results. To do that, firstly we define function
 259 $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ by

260 (4.1)
$$\varphi(a) = \int \Psi(a + w)p(w)dw,$$

261 where $\Psi : \mathbb{R} \rightarrow \mathbb{R}$ is a nonlinear function that satisfies Assumption 3.1, and p is a
 262 probability density function that satisfies Assumptions 2.2 or 2.3.

263 *Remark 4.1.* The mapping φ has all key properties of function Ψ (see Lemma 6.2
 264 in Appendix B in [34], see also [29]). Moreover, it has a strictly positive derivative
 265 at zero, i.e., $\varphi'(0) > 0$, which is necessary to prove our results. The facts that
 266 the nonlinearity Ψ is discontinuous at zero or that it has a positive derivative at
 267 zero, together with condition 4 from Assumptions 2.2 and condition 3 from 2.3, are
 268 crucial to ensure that φ has a positive derivative at zero (see Appendix B in [34], see
 269 also [15, 29]). Notice that the requirement that the pdf p is positive in the vicinity
 270 of the zero is not restrictive, since it holds true for a broad classes of non-zero noise
 271 pdfs.

275 Next, we define functions $\varphi_o : \mathbb{R}^N \rightarrow \mathbb{R}^N$, $\varphi_c : \mathbb{R}^M \rightarrow \mathbb{R}^M$ as $\varphi_o(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N) =$
 276 $[\varphi_o(\mathbf{y}_1), \varphi_o(\mathbf{y}_2), \dots, \varphi_o(\mathbf{y}_N)]$, $\varphi_c(\hat{\mathbf{y}}_1, \hat{\mathbf{y}}_2, \dots, \hat{\mathbf{y}}_M) = [\varphi_c(\hat{\mathbf{y}}_1), \varphi_c(\hat{\mathbf{y}}_2), \dots, \varphi_c(\hat{\mathbf{y}}_M)]$, where
 277 $\mathbf{y} \in \mathbb{R}^N$, $\hat{\mathbf{y}} \in \mathbb{R}^M$ and functions φ_o and φ_c are transformations defined by (4.1)
 278 that correspond to Ψ_o and Ψ_c , respectively. For the a.s. convergence and asymptotic
 279 normality results, we will follow the stochastic approximation framework from [28, 18]
 280 (see Theorem 4 in Appendix A in [34]). That is, we represent algorithm (3.1) in
 281 the form suitable for stochastic approximation analysis. We start by substituting
 282 regression model (2.2) into algorithm (3.2), we get

$$(4.2) \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \alpha_t \left(\frac{b}{a} \mathbf{L}_{\Psi_c}(\mathbf{x}) - \mathbf{H}^\top \Psi_o(\mathbf{H}(\mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \mathbf{n}^t - \mathbf{H}\mathbf{x}^t) \right).$$

283 Define $\boldsymbol{\zeta}^t \in \mathbb{R}^N$ and $\boldsymbol{\eta}^t \in \mathbb{R}^{MN}$ by

$$(4.3)$$

$$286 \quad \boldsymbol{\zeta}^t = \Psi_o(\mathbf{H}(\mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \mathbf{n}^t - \mathbf{H}\mathbf{x}^t) - \varphi_o(\mathbf{H}((\mathbf{1}_N \otimes \boldsymbol{\theta}^*) - \mathbf{x}^t)), \quad \boldsymbol{\eta}^t = \begin{bmatrix} \vdots \\ \sum_{j \in \Omega_i} \boldsymbol{\eta}_{ij}^t \\ \vdots \end{bmatrix},$$

287

288 where $\boldsymbol{\eta}_{ij}^t = \Psi_c(\mathbf{x}_i^t - \mathbf{x}_j^t + \boldsymbol{\xi}_{ij}^t) - \varphi_c(\mathbf{x}_i^t - \mathbf{x}_j^t)$. Now, since φ is defined by (4.1), it can
 289 be shown that $\mathbb{E}[\boldsymbol{\zeta}^t] = \mathbb{E}[\boldsymbol{\eta}^t] = 0$, where the expectation is taken with respect to \mathcal{F}
 290 (see Appendix B in [34]). Furthermore, we define function $\mathbf{L}_{\varphi_c} : \mathbb{R}^{MN} \rightarrow \mathbb{R}^{MN}$ as
 291 $\mathbf{L}_{\varphi_c}(\cdot) = \mathbf{L}_{\Psi_c}(\cdot) - \boldsymbol{\eta}^t$, i.e., its i -th block of size M is $\sum_{j \in \Omega_i} \varphi_c(\mathbf{x}_i - \mathbf{x}_j)$. for $i = 1, 2, \dots, N$.

292 Finally, substituting (4.3) into (4.2), we rewrite algorithm (3.2) by

$$(4.4) \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \alpha_t \left(\frac{b}{a} \mathbf{L}_{\varphi_c}(\mathbf{x}^t) - \mathbf{H}^\top \varphi_o(\mathbf{H}((\mathbf{1}_N \otimes \boldsymbol{\theta}^*) - \mathbf{x}^t)) - \mathbf{H}^\top \boldsymbol{\zeta}^t + \frac{b}{a} \boldsymbol{\eta}^t \right).$$

294

295 Now, we are ready to establish following results.

296 **4.2. Almost sure convergence.** We have the following Theorem.

297 **THEOREM 4.2** (Almost sure convergence). *Let Assumptions 2.1-3.1 hold and*
 298 $\alpha_t = a/(t+1)^\delta$, $\delta \in (0.5, 1]$. *Then, for each agent $i = 1, \dots, N$, the sequence of*
 299 *iterates $\{\mathbf{x}_i^t\}$ generated by algorithm (3.1) converges almost surely to the true vector*
 300 *parameter $\boldsymbol{\theta}^*$.*

301 Theorem 4.2 establishes almost sure convergence of the proposed algorithm (3.1),
 302 whether observation or communication noises have finite or infinite moments of order
 303 greater than one. On the other hand, if we set at least one of the functions Ψ_o, Ψ_c to
 304 be identity functions (and thus recover either the \mathcal{LU} scheme from [18] or the method
 305 from [15]), the resulting method fails to converge (See Appendix D in [34]). In other
 306 words, the methods in [18] and [15] fail to converge under the simultaneous presence
 307 of heavy-tailed observation and communication noises.

308 *Proof.* (Proof of Theorem 4.2)

309 The proof consists of verifying conditions B1–B5 of Theorem 4 in [34] (See Appendix
 310 A in [34]). First, we define quantities $\mathbf{r}(\mathbf{x})$ and $\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)$ by:

$$311 \quad (4.5) \quad \mathbf{r}(\mathbf{x}) = -\frac{b}{a} \mathbf{L}_{\varphi_c}(\mathbf{x}) - \mathbf{H}^\top \varphi_o(\mathbf{H}(\mathbf{x} - (\mathbf{1}_N \otimes \boldsymbol{\theta}^*))),$$

$$312 \quad (4.6) \quad \boldsymbol{\gamma}(t+1, \mathbf{x}, \omega) = -\frac{b}{a} \boldsymbol{\eta}^t + \mathbf{H}^\top \boldsymbol{\zeta}^t.$$

313

314 Here, ω denotes a canonical element of the underlying probability space $(\Omega, \mathcal{F}, \mathbb{P})$.

315 Condition B1 holds because $\mathbf{r}(\cdot)$ is \mathcal{B}^{MN} measurable and $\boldsymbol{\gamma}(t+1, \cdot, \cdot)$ is $\mathcal{B}^{MN} \otimes \mathcal{F}$
 316 measurable for each t , where \mathcal{B}^{MN} is the Borel sigma algebra on \mathbb{R}^{MN} . Consider
 317 the filtration \mathcal{F}_t , $t = 1, 2, \dots$, where \mathcal{F}_t is the σ -algebra generated by $\{\mathbf{n}^s\}_{s=0}^{t-1}$ and

318 $\{\boldsymbol{\xi}_{ij}^s\}_{s=0}^{t-1}$. We have that the family of random vectors $\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)$ is \mathcal{F}_t measurable,
 319 zero-mean and independent of \mathcal{F}_{t-1} . Hence, condition B2 holds.

320 We now show that condition B3 also holds. We use the following Lyapunov function
 321 $V : \mathbb{R}^{MN} \rightarrow \mathbb{R}$,

$$322 \quad (4.7) \quad V(\mathbf{x}) = \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2,$$

324 which is clearly twice continuously differentiable and has uniformly bounded second
 325 order partial derivatives. The gradient of V equals $\nabla V(\mathbf{x}) = 2(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)$. We
 326 must show that

$$327 \quad (4.8) \quad \sup_{\mathbf{x} \in S_\epsilon} \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle < 0,$$

328 where $S_\epsilon = \{\mathbf{x} \in \mathbb{R}^{MN} : \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\| \in (\epsilon, 1/\epsilon)\}$. For any $\mathbf{x} \in \mathbb{R}^{MN}$, we have:

$$329 \quad \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle = 2(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \left(-\frac{b}{a} \mathbf{L}_{\varphi_c}(\mathbf{x}) - \mathbf{H}^\top \varphi_o(\mathbf{H}(\mathbf{x}^t - (\mathbf{1}_N \otimes \boldsymbol{\theta}^*))) \right)$$

$$330 \quad (4.9)$$

$$331 \quad = -\frac{2b}{a} \underbrace{(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{L}_{\varphi_c}(\mathbf{x})}_{T_1(\mathbf{x})} - \underbrace{(\mathbf{H}(\mathbf{x} - (\mathbf{1}_N \otimes \boldsymbol{\theta}^*)))^\top \varphi_o(\mathbf{H}(\mathbf{x} - (\mathbf{1}_N \otimes \boldsymbol{\theta}^*)))}_{T_2(\mathbf{x})}.$$

332 The terms $T_1(\mathbf{x})$ and $T_2(\mathbf{x})$ can be written respectively as

$$333 \quad T_1(\mathbf{x}) = \sum_{\{i,j\} \in E, i < j} (\mathbf{x}_i - \mathbf{x}_j)^\top \varphi_c(\mathbf{x}_i - \mathbf{x}_j) = \sum_{\{i,j\} \in E, i < j} \mathbf{g}^\top \varphi_c(\mathbf{g}),$$

$$334 \quad T_2(\mathbf{x}) = \sum_{i=1}^N \hat{\mathbf{g}}_i^\top \varphi_o(\hat{\mathbf{g}}_i),$$

335 where $\hat{\mathbf{g}} = \mathbf{H}^\top \varphi_o(\mathbf{H}(\mathbf{x}^t - (\mathbf{1}_N \otimes \boldsymbol{\theta}^*)))$, $\mathbf{g} = \mathbf{x}_i - \mathbf{x}_j$ and $\mathbf{g}^\top \varphi_c(\mathbf{g}) = \sum_{\ell=1}^M \mathbf{g}_\ell^\top \varphi_c(\mathbf{g}_\ell)$.
 336 Using the fact that both of the functions φ_c and φ_o are odd functions, for which we
 337 have that $\varphi(a) > 0$ if $a > 0$, we have that $\langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle \geq 0$ for all $\mathbf{x} \in \mathbb{R}^{MN}$ (see
 338 Appendix B in [34]). Moreover, recalling the fact that function φ_c is continuous at
 339 zero, and equal to zero only at zero, we have that $T_1(\mathbf{x})$ is equal to zero if and only if
 340 $\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^* = \mathbf{1}_N \otimes \mathbf{m}$, for $\mathbf{m} \in \mathbb{R}^M$ (see Lemma 6 in Appendix B in [34]). We only
 341 consider the case when $\mathbf{m} \neq 0$, since from $\mathbf{m} = 0$ we have that $\mathbf{x} = \mathbf{1}_N \otimes \boldsymbol{\theta}^*$, which is
 342 not in the set S_ϵ . However, for that choice of $\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*$ we have that

$$343 \quad T_2(\mathbf{1}_N \otimes \boldsymbol{\theta}^* + \mathbf{1}_N \otimes \mathbf{m}) = (\mathbf{H} \mathbf{1}_N \otimes \mathbf{m})^\top \varphi_o(\mathbf{H} \mathbf{1}_N \otimes \mathbf{m})$$

$$344 \quad = \sum_{i=1}^N (\mathbf{h}_i^\top \mathbf{m}) \varphi_o(\mathbf{h}_i^\top \mathbf{m}) > 0,$$

345 since $\mathbf{h}_i^\top \mathbf{m}$ and $\varphi_o(\mathbf{h}_i^\top \mathbf{m})$ have the same sign. Hence, for all $\epsilon > 0$ we have that
 346 $\sup_{\mathbf{x} \in S_\epsilon} \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle < 0$. Thus, condition B3 also holds.

347 Now we inspect condition B4. From equation (4.5) we have that

$$348 \quad (4.10) \quad \|\mathbf{r}(\mathbf{x})\|^2 \leq \left\| \frac{b}{a} \mathbf{L}_{\varphi_c}(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) \right\|^2 + \|\mathbf{H}^\top \varphi_o(\mathbf{H}(\mathbf{x} - (\mathbf{1}_N \otimes \boldsymbol{\theta}^*)))\|^2 \leq c_1(1 + V(\mathbf{x})),$$

349 for some positive constant c_1 (see Appendix B in [34]). Moreover, we have that

$$350 \quad (4.11) \quad \|\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)\|^2 \leq \left\| \frac{b}{a} \boldsymbol{\eta}^t \right\|^2 + \|\mathbf{H}^\top \boldsymbol{\zeta}^t\|^2$$

351 which leads to

$$352 \quad (4.12) \quad \mathbb{E} \left[\|\boldsymbol{\gamma}(t+1, \mathbf{x}^t, \omega)\|^2 \right] \leq c_2(1 + V(\mathbf{x})),$$

353

359 for some positive constant c_2 . Finally, we have that

$$360 \quad \|\mathbf{r}(\mathbf{x})\|^2 + \mathbb{E} \left[\|\gamma(t+1, \mathbf{x}^t, \omega)\|^2 \right] \leq c_3(1 + V(\mathbf{x})),$$

362 for some positive constant c_3 . Setting that $\epsilon \rightarrow 0^+$ in (4.8), for all $\mathbf{x} \in \mathbb{R}^{MN}$, we have
363 that $\langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle \leq 0$. Thus,

$$364 \quad \|\mathbf{r}(\mathbf{x})\|^2 + \mathbb{E} \left[\|\gamma(t+1, \mathbf{x}^t, \omega)\|^2 \right] \leq c_3(1 + V(\mathbf{x})) - k \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle$$

366 for every $k > 0$. Therefore, condition B4 also holds. Condition B5 holds by the
367 definition of the algorithm (3.1). Thus, almost sure convergence is proved. \square

368 **4.3. Asymptotic normality.** We now consider asymptotic normality of the
369 proposed estimator (3.1). We have the following theorem.

370 **THEOREM 4.3** (Asymptotic normality). *Let Assumptions 2.1-3.1 hold. Consider*
371 *algorithm (3.1) with step-size $\alpha_t = a/(t+1)^\delta$, $t = 0, 1, \dots$, $a > 0$, with $\delta = 1$. Then,*
372 *the normalized sequence of iterates $\{\sqrt{t+1}(\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)\}$ converges in distribution*
373 *to a zero-mean multivariate normal random vector, i.e., the following holds:*

$$374 \quad \sqrt{t+1}(\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) \Rightarrow \mathcal{N}(\mathbf{0}, \mathbf{S}),$$

376 where the asymptotic covariance matrix \mathbf{S} equals:

$$377 \quad (4.13) \quad \mathbf{S} = a^2 \int_0^\infty e^{\Sigma v} \mathbf{S}_0 e^{\Sigma^\top v} dv.$$

378 Here, $\mathbf{S}_0 = \frac{b^2}{a^2} \sigma_c^2 \text{Diag}(\{d_i \mathbf{I}_M\}) - \frac{b}{a} \mathbf{K}_{c,o} \mathbf{H} - \frac{b}{a} \mathbf{H}^\top \mathbf{K}_{c,o}^\top + \sigma_o^2 \mathbf{H}^\top \mathbf{H}$; $\sigma_o^2 = \int |\Psi_o(w)|^2$
379 $d\Phi_o(w)$ is the effective observation noise variance after passing through the nonlin-
380 earity Ψ_o ; $\sigma_c^2 = \int |\Psi_c(w)|^2 d\Phi_c(w)$ is the effective communication noise variance after
381 passing through the nonlinearity Ψ_c ; $\mathbf{K}_{c,o} \in \mathbf{R}^{MN \times N}$ is the effective cross-covariance
382 matrix between the observation and the communication noise after passing through
383 the appropriate nonlinearity, i.e., the (k, s) element of the matrix $\mathbf{K}_{c,o}$ is given by
384 $[(\mathbf{K}_{c,o})]_{ks} = \sum_{j \in \Omega_i} \int \int \Psi_c(w_{ij\ell}) \Psi_o(w_k) p_{k,ij\ell}^{c,o}(w_{ij\ell}, w_k) dw_{ij\ell} dw_k$. Here, ℓ satisfies the fol-

385 lowing: $s = M(i-1) + \ell$; and $p_{k,ij\ell}^{c,o}$ is the joint probability density function for the k -th
386 observation noise n_k and the ℓ -th element of the communication noise $[(\boldsymbol{\xi}_{ij})]_\ell$. We
387 also recall the observation matrix \mathbf{H} in (2.2); functions φ_c, φ_o appropriate versions
388 of function φ in (4.1); and $\Sigma = \frac{1}{2} \mathbf{I} - a(\frac{b}{a} \varphi'_c(0) \mathbf{L} \otimes \mathbf{I}_M + \varphi'_o(0) \mathbf{H}^\top \mathbf{H})$; here, a is taken
389 large enough such that matrix Σ is stable.

390 **Remark 4.4.** Notice that, for the assumed setting, σ_c^2 and σ_o^2 are finite. Also,
391 $\mathbf{K}_{c,o}$ is finite, i.e., $\|\mathbf{K}_{c,o}\| < \infty$, since we have that

$$392 \quad \left| \int \int \Psi_c(w_1) \Psi_o(w_2) d\Phi^{c,o} \right| \leq \int \int |\Psi_c(w_1) \Psi_o(w_2)| d\Phi^{c,o} < \frac{1}{2} \sigma_c^2 + \frac{1}{2} \sigma_o^2.$$

394 **Remark 4.5.** If we assume that observation and communication noise are mutu-
395 ally independent, the only difference from the previous theoretical results occurs in
396 the $\mathbf{A}(t, \mathbf{x})$, i.e., in the \mathbf{S}_0 . Under this setting, matrix \mathbf{S}_0 is now equal to

$$397 \quad \mathbf{S}_0 = \frac{b^2}{a^2} \sigma_c^2 \text{Diag}(\{d_i \mathbf{I}_M\}) + \sigma_o^2 \mathbf{H}^\top \mathbf{H},$$

399 which is expected, since the effective cross-covariance matrix $\mathbf{K}_{c,o}$ is now equal to
400 zero.

401 **Theorem 4.3** establishes asymptotic normality of the proposed method. This is
402 achieved with heavy-tailed observation and communication noise an the nonlinearities
403 Ψ_o and Ψ_c with uniformly bounded outputs. Moreover, the theorem explicitly
404 evaluates the corresponding asymptotic variance. When the two noises are mutually
405 independent, Ψ_o is identity, and observation noise variance is finite, we recover the re-

406 sult in [15], Theorem 3.5, as a special case. That is, a notable difference with respect
 407 to [15] is the ability to handle here mutually correlated observation and communi-
 408 cation noises. The effect of correlation is complex in general, however, as shown in
 409 Section 5 later, generally a stronger positive noises correlation leads to a lower asymp-
 410 totic variance. Intuitively, at an extreme, a full positive correlation practically means
 411 that only one effective noise exists in the system, and hence it can be suppressed more
 412 easily. Further, note that Theorem 4.3 establishes a local asymptotic rate $O(1/t)$ of
 413 \mathbf{x}^t to zero, in the weak convergence sense, when $\alpha_t = a/(t+1)$. We show later (see
 414 Theorem 4.6) that a global MSE rate $O(1/t^{\hat{\delta}})$ with a lower (worse) degree $\hat{\delta}$ can be
 415 established when step-size $\alpha_t = a/(t+1)^{\delta}$, $\delta \in (0.5, 1)$, is used.

416 We next discuss asymptotic efficiency² of the proposed estimator. We first briefly
 417 review the relevant existing work to better position our results. First, consider the
 418 best linear centralized estimator $\mathbf{x}_{\text{cent}}^t$ of $\boldsymbol{\theta}^*$, that has access to measurements from all
 419 sensors (nodes) $n = 1, 2, \dots, N$ at all times $t = 0, 1, \dots$. In the general case, addition-
 420 ally assuming that observation noise has finite variance, the best linear centralized
 421 estimator $\mathbf{x}_{\text{cent}}^t$ is asymptotically normal and has the lowest asymptotic covariance
 422 matrix \mathbf{S}_{cent} among all estimators of $\boldsymbol{\theta}^*$ when the only knowledge of observation
 423 noise is variance and no other information of noise distribution is known. Moreover,
 424 its asymptotic covariance matrix \mathbf{S}_{cent} attains the Cramér-Rao lower bound if the
 425 observation noise is Gaussian (see for example [17]). On the other hand, when the
 426 probability density function is known, the centralized estimator in [29] can be tuned to
 427 the pdf of the observation noise so that it achieves the Cramér-Rao bound. In the dis-
 428 tributed setting, when there is no communication noise, the authors of [17] develop an
 429 estimator which is asymptotically normal and has the optimal asymptotic covariance
 430 matrix \mathbf{S}_{cent} (optimal in the sense that the asymptotic covariance matrix is the same
 431 as for the best linear centralized estimator $\mathbf{x}_{\text{cent}}^t$). We now discuss the asymptotic
 432 covariance matrix \mathbf{S} of the proposed estimator (3.1). This quantity depends on the
 433 system parameters, including network topology and communication noise. Therefore,
 434 in the general case, the proposed estimator (3.1) is not asymptotically efficient, i.e.,
 435 $\mathbf{S} \neq \mathbf{I}^{-1}(\boldsymbol{\theta}^*)$, where $\mathbf{I}(\boldsymbol{\theta}^*)$ is the Fisher information matrix. However, with respect to
 436 the proposed distributed recursive estimator, we make the following observations. 1)
 437 First, the estimator is order-optimal in the weak convergence sense; that is, its (weak
 438 convergence sense) rate of error decay is the same as that of the asymptotically ef-
 439 ficient estimator. 2) The corresponding ‘‘convergence constant,’’ i.e., the asymptotic
 440 covariance, is different from that of the centralized Cramér-Rao-optimal estimator,
 441 and it is hence not optimal. We note that the paper provides major contributions
 442 with respect to state of the art, as it gives the first distributed estimator that ensures
 443 almost sure convergence in the presence of infinite variance correlated sensing and
 444 communication noises; moreover, its weak convergence sense rate of convergence is
 445 order-optimal. It remains an interesting future work direction to explore whether an
 446 optimal asymptotic covariance can be achieved in this setting via distributed estima-
 447 tors. In view of the results [29] for the centralized setting, it is likely that this cannot
 448 be achieved unless the nonlinearities are tuned to the noise pdfs that in turn have to
 449 be known.

450 *Proof.* (Proof of Theorem 4.3)

²An estimator \mathbf{y}^t of an unknown parameter $\boldsymbol{\theta}^*$, for which we have that $\sqrt{t+1}(\mathbf{y}^t - \boldsymbol{\theta}^*) \Rightarrow \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$, is said to be asymptotically efficient if $\mathbf{S} = \mathbf{I}^{-1}(\boldsymbol{\theta}^*)$, where $\mathbf{I}(\boldsymbol{\theta}^*)$ is the Fisher information matrix. The Fisher information matrix represents the best achievable asymptotic covariance by any estimator, as determined by the well-known Cramer-Rao bound (see [28]).

451 We prove Theorem 4.3 in the same manner as Theorem 4.2 is proved, i.e., by verifying
 452 assumptions C1-C5 of Theorem 4 in [34] (see Appendix A in [34]). Function $\mathbf{r}(\cdot)$
 453 defined by (4.5) can be written as

$$454 \quad \mathbf{r}(\mathbf{x}) = -\frac{b}{a}\varphi'_c(0)\mathbf{L} \otimes \mathbf{I}_M (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) - \varphi'_o(0)\mathbf{H}^\top \mathbf{H} (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \boldsymbol{\delta}(\mathbf{x}),$$

455 Here, mapping $\boldsymbol{\delta} : \mathbb{R}^{MN} \rightarrow \mathbb{R}^{MN}$ is given by:

$$457 \quad (4.14) \quad \boldsymbol{\delta}(\mathbf{x}) = -\frac{b}{a}\mathbf{L}\boldsymbol{\delta}_c(\mathbf{x}) - \mathbf{H}^\top \boldsymbol{\delta}_o(\mathbf{H}(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)).$$

458 Next, mapping $\mathbf{L}\boldsymbol{\delta}_c(\mathbf{x}) : \mathbb{R}^{MN} \rightarrow \mathbb{R}^{MN}$ is vector of size MN such that the i -th M -size
 459 block equals $\sum_{j \in \Omega_i} \boldsymbol{\delta}_c(\mathbf{x}_i - \mathbf{x}_j)$, $i = 1, 2, \dots, N$, mappings $\boldsymbol{\delta}_c : \mathbb{R}^M \rightarrow \mathbb{R}^M$, $\boldsymbol{\delta}_o : \mathbb{R}^N \rightarrow$

461 \mathbb{R}^N are component-wise maps of $\boldsymbol{\delta}_c$ and $\boldsymbol{\delta}_o$ are first order residuals that corresponds
 462 to φ_c and φ_o respectively, i.e., $\boldsymbol{\delta}_c(\mathbf{y}_1, \mathbf{y}_1, \dots, \mathbf{y}_M) = [\delta_c(\mathbf{y}_1), \delta_c(\mathbf{y}_2), \dots, \delta_c(\mathbf{y}_M)]^\top$ and
 463 $\boldsymbol{\delta}_o(\hat{\mathbf{y}}_1, \hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_N) = [\delta_o(\hat{\mathbf{y}}_1), \delta_o(\hat{\mathbf{y}}_2), \dots, \delta_o(\hat{\mathbf{y}}_M)]^\top$ for $\mathbf{y} \in \mathbb{R}^N$, $\hat{\mathbf{y}} \in \mathbb{R}^M$ (see Appendix
 464 B in [34]).

465 Thus, $\mathbf{r}(\mathbf{x})$ admits representation in (36) of Theorem 4 in [34] for $\mathbf{B} = -\frac{b}{a}\varphi'_c(0)\mathbf{L} \otimes$
 466 $\mathbf{I}_M - \varphi'_o(0)\mathbf{H}^\top \mathbf{H}$ and mapping $\boldsymbol{\delta}(\cdot)$ defined by (4.14). Therefore, condition C1 holds.
 467 Since we use that $\alpha_t = \frac{a}{t+1}$, condition C2 trivially holds. Furthermore, $\boldsymbol{\Sigma} = a\mathbf{B} + \frac{1}{2}\mathbf{I}$
 468 is stable if a is large enough, because matrix $-\mathbf{B}$ is positive definite (See [18]). Thus,
 469 condition C3 also holds.

470 For $\mathbf{A}(t, \mathbf{x}) = \mathbb{E}[\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)\boldsymbol{\gamma}^\top(t+1, \mathbf{x}, \omega)]$ it is easy to show that

$$471 \quad \lim_{t \rightarrow \infty, \mathbf{x} \rightarrow \boldsymbol{\theta}^*} \mathbf{A}(t, \mathbf{x}) = \frac{b^2}{a^2}\sigma_c^2 \text{Diag}(\{d_i \mathbf{I}_M\}) - \frac{b}{a}\mathbf{K}_{c,o}\mathbf{H} - \frac{b}{a}\mathbf{H}^\top \mathbf{K}_{c,o}^\top + \sigma_o^2 \mathbf{H}^\top \mathbf{H}.$$

472 Therefore, condition C4 also holds. To show that condition C5 holds, it is suffice
 473 to show that the family of random variables $\{\|\boldsymbol{\gamma}_\varphi(t+1, \mathbf{x}, \omega)\|^2\}_{t=0,1,\dots, \|\mathbf{x}-\boldsymbol{\theta}^*\|<\epsilon}$ is
 474 uniformly integrable. To do that, follow the arguments as in e.g., [18] and [15]. \square
 475

476 **4.4. Mean squared error convergence.** In this subsection, we state and prove
 477 a result on the mean squared error (MSE) convergence rate when both nonlinearities
 478 Ψ_o and Ψ_c satisfy part 5' of Assumption 3.1, i.e., $|\Psi_o| \leq c_o$, $|\Psi_c| \leq c_c$, for some positive
 479 constants c_o and c_c . Moreover, we set the step size to $\alpha_t = \frac{a}{(t+1)^\delta}$, for $\delta \in (\frac{1}{2}, 1)$. We
 480 have the following theorem.

481 **THEOREM 4.6 (MSE convergence).** *Let Assumptions 2.1-3.1 hold. Then, for*
 482 *the sequence of iterates $\{\mathbf{x}^t\}$ generated by algorithm (3.2), provided that the step-size*
 483 *sequence $\{\alpha_t\}$ is given by $\alpha_t = a/(t+1)^\delta$, $a > 0, \delta \in (0.5, 1)$, there exists $\hat{\delta} \in (0, 1)$*
 484 *such that $\mathbb{E}[\|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2] = O(1/t^{\hat{\delta}})$.*

485 Theorem 4.6 establishes a MSE convergence rate of the proposed estimator (3.2) under
 486 the simultaneous presence of heavy-tailed (possibly infinite variance) observation and
 487 communication noises, when both the observation and communication nonlinearities
 488 have uniformly bounded outputs. This is in contrast with recent studies on distributed
 489 estimation in heavy-tailed noises like [15] that only establishes a.s. and asymptotic
 490 normality results. We refer to the proof of Theorem 4.6 for the exact value of the
 491 convergence rate power $\hat{\delta}$.

492 **Setting up the proof.** We now prove Theorem 4.6 through a sequence of
 493 intermediate results (Lemmas). Recall quantities $\mathbf{r}(\cdot)$, $\boldsymbol{\gamma}(\cdot, \cdot, \cdot)$ and $V(\cdot)$ from (4.5),
 494 (4.6) and (4.7) respectively. The proof will be based on establishing a sufficient decay
 495 on quantity $\mathbb{E}[V(\mathbf{x}^t)]$. First, notice that algorithm (4.4) can be written as

$$496 \quad \mathbf{x}^{t+1} = \mathbf{x}^t + \alpha_t (\mathbf{r}(\mathbf{x}^t) + \boldsymbol{\gamma}(t+1, \mathbf{x}^t, \omega)).$$

498 Moreover, we have that

$$499 \quad V(\mathbf{x}^{t+1}) = V(\mathbf{x}^t) + 2\alpha_t (\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top (\mathbf{r}(\mathbf{x}^t) + \boldsymbol{\gamma}(t+1, \mathbf{x}^t, \omega)) \\ 500 \quad \quad \quad + \alpha_t^2 \|\mathbf{r}(\mathbf{x}^t) + \boldsymbol{\gamma}(t+1, \mathbf{x}^t, \omega)\|^2 \\ 501 \quad \quad \quad = V(\mathbf{x}^t) + 2\alpha_t (\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top (\mathbf{r}(\mathbf{x}^t) + \boldsymbol{\gamma}(t+1, \mathbf{x}^t, \omega)) + \alpha_t^2 c',$$

502 for positive constant $c' = \|\mathbf{r}(\mathbf{x}^t) + \boldsymbol{\gamma}(t+1, \mathbf{x}^t, \omega)\|^2 < \infty$. Therefore, taking a condi-
503 tional expectation with respect to \mathcal{F}_t , we have:

$$504 \quad (4.15) \quad \mathbb{E}[V(\mathbf{x}^{t+1})|\mathcal{F}_t] = V(\mathbf{x}^t) + 2\alpha_t (\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{r}(\mathbf{x}^t) + \alpha_t^2 c'.$$

505 Also, from equation (4.9), it follows that

$$506 \quad (4.16) \quad (\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{r}(\mathbf{x}^t) = -\frac{b}{a} T_1(\mathbf{x}^t) - T_2(\mathbf{x}^t).$$

507 We next need to show that the quantity in (4.16) is ‘‘sufficiently negative’’, relative
508 to quantity $V(\mathbf{x}^t)$. This is achieved through a sequence of lemmas. First, we upper
509 bound quantities $\|\mathbf{x}^t\|$ and $\|\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|$.

510 LEMMA 4.7. *Let Assumptions 2.1-3.1 hold. Then, for the sequence of iterates*
511 *$\{\mathbf{x}^t\}$ generated by algorithm (3.2), provided that the step-size sequence $\{\alpha_t\}$ is given*
512 *by $\alpha_t = a/(t+1)^\delta$, $a > 0, \delta \in (0.5, 1)$, we have that, for any outcome ω :*

$$513 \quad (4.17) \quad \|\mathbf{x}^t\| \leq g_t = \|\mathbf{x}^0\| + \left(b\sqrt{MN}dc_c + a\|\mathbf{H}\|\sqrt{N}c_o \right) \frac{t^{1-\delta}}{1-\delta},$$

$$(4.18)$$

$$514 \quad \|\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\| \leq g'_t = \|\mathbf{x}^0 - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\| + \left(b\sqrt{MN}dc_c + a\|\mathbf{H}\|\sqrt{N}c_o \right) \frac{t^{1-\delta}}{1-\delta}.$$

515 Consequently, $\|\mathbf{H}(\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)\| \leq \|\mathbf{H}\|g'_t$.

516 *Proof.* Using the boundness of the nonlinearities, we have that $\|\mathbf{L}_{\Psi_c}(\mathbf{x})\|^2 \leq$
517 $\sqrt{MN}dc_c$ and $\|\mathbf{H}^\top \Psi_o(\mathbf{H}(\mathbf{1}_N \otimes \boldsymbol{\theta}^* - \mathbf{x}^t) + \mathbf{n}^t)\| \leq \|\mathbf{H}\|\sqrt{N}c_o$, where $d = \max_i d_i$.

518 Therefore, recalling the algorithm (4.2), for all $t > 0$ we have that

$$519 \quad \|\mathbf{x}^t\| \leq \|\mathbf{x}^{t-1}\| + \underbrace{\alpha_{t-1} \left(\frac{b}{a}\sqrt{MN}dc_c + \|\mathbf{H}\|\sqrt{N}c_o \right)}_c \leq \|\mathbf{x}^{t-2}\| + \alpha_{t-2}c + \alpha_{t-1}c$$

$$520 \quad \leq \|\mathbf{x}^0\| + c \sum_{j=0}^{t-1} \frac{a}{(1+j)^\delta} \leq \|\mathbf{x}^0\| + c \int_0^{t-1} \frac{a}{(1+s)^\delta} ds \leq \|\mathbf{x}^0\| + ca \frac{t^{1-\delta}}{1-\delta}.$$

521 Analogously, for all $t > 0$, we have that $\|\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\| \leq g'_t$, and as a consequence
522 $\|\mathbf{H}(\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)\| \leq \|\mathbf{H}\|g'_t$. \square

523 Next, we have the following Lemma that bounds quantities $T_1(x)$ and $T_2(x)$.

524 LEMMA 4.8. *Let Assumptions 2.1-3.1 hold. Then, for the sequence of iterates*
525 *$\{\mathbf{x}^t\}$ generated by algorithm (3.2), provided that the step-size sequence $\{\alpha_t\}$ is given*
526 *by $\alpha_t = a/(t+1)^\delta$, $a > 0, \delta \in (0.5, 1)$, we have that there exist positive constants G_c*
527 *and G_o such that, for any outcome ω :*

$$528 \quad T_1(\mathbf{x}^t) \geq \frac{\varphi'_c(0)G_c}{4g_t} (\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{L} \otimes \mathbf{I} (\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*),$$

$$529 \quad T_2(\mathbf{x}^t) \geq \frac{\varphi'_o(0)G_o}{2\|\mathbf{H}\|g'_t} (\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{H}^\top \mathbf{H} (\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*),$$

530 To prove Lemma 4.8, we make use of the following Lemma from [13] (see Lemma 5.5
531 in [13]).

532 LEMMA 4.9. *Consider function φ in (4.1), there exists a positive constant G such*

539 that $|\varphi(a)| \leq \frac{\varphi'(0)G|a|}{2g}$, for all $|a| < g$.

540 *Proof.* (Proof of Theorem 4.8) Using Lemma 4.9 for function φ_c we get that there
 541 exists a positive constant G_c such that

$$\begin{aligned}
 542 \quad T_1(\mathbf{x}^t) &= \sum_{\{i,j\} \in E, i < j} (\mathbf{x}_i^t - \mathbf{x}_j^t)^\top \varphi_c(\mathbf{x}_i^t - \mathbf{x}_j^t) \\
 543 \quad &= \sum_{\{i,j\} \in E, i < j} \sum_{\ell=1}^M ((\mathbf{x}_i^t)_\ell - (\mathbf{x}_j^t)_\ell)^\top \varphi_c((\mathbf{x}_i^t)_\ell - (\mathbf{x}_j^t)_\ell) \\
 544 \quad (4.19) \quad &\geq \frac{\varphi'_c(0)G_c}{4g_t} \sum_{\{i,j\} \in E, i < j} \|\mathbf{x}_i^t - \mathbf{x}_j^t\|^2 = \frac{\varphi'_c(0)G_c}{4g_t} (\mathbf{x}^t)^\top (\mathbf{L} \otimes \mathbf{I}) \mathbf{x}^t \\
 545 \quad &= \frac{\varphi'_c(0)G_c}{4g_t} (\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{L} \otimes \mathbf{I} (\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*),
 \end{aligned}$$

546 since, from Lemma 4.7 we have $\|\mathbf{x}^t\| \leq g_t$. Analogously, from Lemma 4.9 we have
 547 that for the function φ_o there exists a positive constant G_o such that

$$\begin{aligned}
 549 \quad T_2(\mathbf{x}) &= \sum_{i=1}^N (\mathbf{H}(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*))_i \varphi_o((\mathbf{H}(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*))_i) \\
 550 \quad (4.20) \quad &\geq \frac{\varphi'_o(0)G_o}{2\|\mathbf{H}\|g'_t} (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{H}^\top \mathbf{H} (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*),
 \end{aligned}$$

551 since, from Lemma 4.7 we have $\|\mathbf{H}(\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)\| \leq \|\mathbf{H}\|g'_t$. \square

552 We next have the following theorem that analyzes positive definiteness of the
 553 matrix $\frac{\varphi'_c(0)G_c}{4g_t} \mathbf{L} \otimes \mathbf{I} + \frac{\varphi'_o(0)G_o}{2\|\mathbf{H}\|g'_t} \mathbf{H}^\top \mathbf{H}$.

554 LEMMA 4.10. *Let Assumptions 2.1-3.1 hold. The following is true for any $\mathbf{x} \in$*
 555 \mathbb{R}^{MN} :

$$\begin{aligned}
 557 \quad &(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \left(\frac{\varphi'_c(0)G_c}{4g_t} \mathbf{L} \otimes \mathbf{I} + \frac{\varphi'_o(0)G_o}{2\|\mathbf{H}\|g'_t} \mathbf{H}^\top \mathbf{H} \right) (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) \\
 558 \quad &\geq \min \left\{ \frac{\varphi'_o(0)G_o}{2\|\mathbf{H}\|g'_t} \left(\frac{\lambda_H}{N} - \frac{2S_H}{\sqrt{N}}k \right), \frac{b\varphi'_c(0)G_c}{4ag_t} \frac{\lambda_2(\mathbf{L})}{1 + \frac{1}{k^2}} \right\} \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2,
 \end{aligned}$$

560 where g_t and g'_t are defined in Lemma 4.7, G_c and G_o in Lemma 4.8, $S_H = \sum_{i=1}^N \|\mathbf{h}_i\|^2$,

561 $\lambda_H = \lambda_1 \left(\sum_{i=1}^N \mathbf{h}_i \mathbf{h}_i^\top \right) > 0$ is the smallest eigenvalue of regular matrix $\sum_{i=1}^N \mathbf{h}_i \mathbf{h}_i^\top$ (see
 562 Assumption 2.1) and recalling that $\lambda_2(\mathbf{L}) > 0$ is the smallest positive eigenvalue of
 563 Laplacian matrix \mathbf{L} .

564 *Proof.* Let us consider matrix $\mathbf{L} \otimes \mathbf{I} + \mathbf{H}^\top \mathbf{H}$ and follow argument as in Appendix A
 565 of [32]. For any $\mathbf{x} \in \mathbb{R}^{MN}$, we have that there exist vectors $\mathbf{u} \in \text{span}\{\mathbf{1} \otimes \mathbf{m} | \mathbf{m} \in \mathbb{R}^M\}$
 566 and $\mathbf{v} \in \text{span}\{\mathbf{1} \otimes \mathbf{m} | \mathbf{m} \in \mathbb{R}^M\}^\perp$ such that $\mathbf{x} = \mathbf{u} + \mathbf{v}$. Firstly, we have that

$$\begin{aligned}
 567 \quad &(\mathbf{u} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{H}^\top \mathbf{H} (\mathbf{u} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) = \sum_{i=1}^N (\hat{\mathbf{u}} - \boldsymbol{\theta}^*)^\top \mathbf{h}_i \mathbf{h}_i^\top (\hat{\mathbf{u}} - \boldsymbol{\theta}^*) \\
 568 \quad &= (\hat{\mathbf{u}} - \boldsymbol{\theta}^*)^\top \left(\sum_{i=1}^N \mathbf{h}_i \mathbf{h}_i^\top \right) (\hat{\mathbf{u}} - \boldsymbol{\theta}^*) \\
 569 \quad &\geq \lambda_H \|\hat{\mathbf{u}} - \boldsymbol{\theta}^*\|^2,
 \end{aligned}$$

570 where $\hat{\mathbf{u}} \in \mathbb{R}^M$ such that $\mathbf{u} = \mathbf{1} \otimes \hat{\mathbf{u}}$. Notice here that $\|\mathbf{u} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\| = \sqrt{N} \|\hat{\mathbf{u}} - \boldsymbol{\theta}^*\|$.
 571 Secondly, $(\mathbf{x} - \mathbf{u})^\top \mathbf{H}^\top \mathbf{H} (\mathbf{x} - \mathbf{u}) \geq 0$, since $\mathbf{H}^\top \mathbf{H}$ is positive semi-definite matrix.

573 Thirdly, following also holds

$$\begin{aligned}
574 \quad (\mathbf{u} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{H}^\top \mathbf{H} (\mathbf{x} - \mathbf{u}) &= \sum_{i=1}^N (\hat{\mathbf{u}} - \boldsymbol{\theta}^*)^\top \mathbf{h}_i \mathbf{h}_i^\top (\mathbf{x}_i - \hat{\mathbf{u}}) \\
575 &\geq - \sum_{i=1}^N \|\hat{\mathbf{u}} - \boldsymbol{\theta}^*\| \|\mathbf{h}_i\|^2 \|\mathbf{x}_i - \hat{\mathbf{u}}\| \\
576 &\geq - \|\hat{\mathbf{u}} - \boldsymbol{\theta}^*\| \|\mathbf{v}\| S_{\mathbf{H}}.
\end{aligned}$$

578 Analogously, we have that $(\mathbf{x} - \mathbf{u})^\top \mathbf{H}^\top \mathbf{H} (\mathbf{u} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) \geq -\|\hat{\mathbf{u}} - \boldsymbol{\theta}^*\| \|\mathbf{v}\| S_{\mathbf{H}}$. There-
579 fore,

$$580 \quad (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{H}^\top \mathbf{H} (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) \geq \lambda_{\mathbf{H}} \|\hat{\mathbf{u}} - \boldsymbol{\theta}^*\|^2 - 2S_{\mathbf{H}} \|\hat{\mathbf{u}} - \boldsymbol{\theta}^*\| \|\mathbf{v}\|.$$

582 We also have that $\mathbf{u} - \mathbf{1} \otimes \boldsymbol{\theta}^* \in \text{null}(\mathbf{L} \otimes \mathbf{I})$ and $\mathbf{v} \in \text{Range}(\mathbf{L} \otimes \mathbf{I})$ and, hence, we have
583 that

$$\begin{aligned}
584 \quad (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{L} \otimes \mathbf{I} (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) &= (\mathbf{u} - \mathbf{1}_N \otimes \boldsymbol{\theta}^* + \mathbf{v})^\top \mathbf{L} \otimes \mathbf{I} (\mathbf{u} - \mathbf{1}_N \otimes \boldsymbol{\theta}^* + \mathbf{v}) \\
585 &= \mathbf{v}^\top \mathbf{L} \otimes \mathbf{I} \mathbf{v} \geq \lambda_2(\mathbf{L} \otimes \mathbf{I}) \|\mathbf{v}\|^2 = \lambda_2(\mathbf{L}) \|\mathbf{v}\|^2.
\end{aligned}$$

587 Let $k > 0$ be arbitrarily chosen. If $\|\mathbf{v}\| \leq k \|\mathbf{u} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|$, then we have that

$$\begin{aligned}
588 \quad (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top (\mathbf{L} \otimes \mathbf{I} + \mathbf{H}^\top \mathbf{H}) (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) \\
589 &\geq \lambda_{\mathbf{H}} \|\hat{\mathbf{u}} - \boldsymbol{\theta}^*\|^2 - 2S_{\mathbf{H}} \|\hat{\mathbf{u}} - \boldsymbol{\theta}^*\| \|\mathbf{v}\| + \lambda_2(\mathbf{L}) \|\mathbf{v}\|^2 \\
590 &\geq \left(\frac{\lambda_{\mathbf{H}}}{N} - \frac{2S_{\mathbf{H}}}{\sqrt{N}} k \right) \|\mathbf{u} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2 + \lambda_2(\mathbf{L}) \|\mathbf{v}\|^2 \\
591 &\geq \min \left\{ \frac{\lambda_{\mathbf{H}}}{N} - \frac{2S_{\mathbf{H}}}{\sqrt{N}} k, \lambda_2(\mathbf{L}) \right\} \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2,
\end{aligned}$$

593 where in the last inequality we used the fact that $\|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2 = \|\mathbf{u} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2 +$
594 $\|\mathbf{v}\|^2$. If $\|\mathbf{v}\| \geq k \|\mathbf{u} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|$, then

$$\begin{aligned}
595 \quad (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top (\mathbf{L} \otimes \mathbf{I} + \mathbf{H}^\top \mathbf{H}) (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) &\geq 0 + \lambda_2(\mathbf{L}) \|\mathbf{v}\|^2 \\
596 &\geq \frac{\lambda_2(\mathbf{L})}{1 + \frac{1}{k^2}} \|\mathbf{v}\|^2 + \frac{\lambda_2(\mathbf{L})}{1 + \frac{1}{k^2}} \|\mathbf{u} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2 \\
597 &\geq \frac{\lambda_2(\mathbf{L})}{1 + \frac{1}{k^2}} \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2.
\end{aligned}$$

598 Therefore, regardless of vector \mathbf{v} , we have that

$$\begin{aligned}
600 \quad (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top (\mathbf{L} \otimes \mathbf{I} + \mathbf{H}^\top \mathbf{H}) (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) \\
601 &\geq \min \left\{ \frac{\lambda_{\mathbf{H}}}{N} - \frac{2S_{\mathbf{H}}}{\sqrt{N}} k, \frac{\lambda_2(\mathbf{L})}{1 + \frac{1}{k^2}} \right\} \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2.
\end{aligned}$$

602 Following the same idea, we get that

$$\begin{aligned}
604 \quad (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \left(\frac{\varphi'_c(0)G_c}{4g_t} \mathbf{L} \otimes \mathbf{I} + \frac{\varphi'_o(0)G_o}{2\|\mathbf{H}\|g'_t} \mathbf{H}^\top \mathbf{H} \right) (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) \\
605 \quad (4.21) \quad \geq \min \left\{ \frac{\varphi'_o(0)G_o}{2\|\mathbf{H}\|g'_t} \left(\frac{\lambda_{\mathbf{H}}}{N} - \frac{2S_{\mathbf{H}}}{\sqrt{N}} k \right), \frac{b \varphi'_c(0)G_c}{4ag_t} \frac{\lambda_2(\mathbf{L})}{1 + \frac{1}{k^2}} \right\} \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2. \quad \square
\end{aligned}$$

607 Finally, to prove Theorem 4.6, we make use of the following Lemma from [13] (see
608 Theorem 5.2 in [13]).

609 LEMMA 4.11. *Let z^t be a nonnegative (deterministic) sequence satisfying*

$$610 \quad z^{t+1} \leq (1 - r_1^t) z^t + r_2^t,$$

612 *for all $t \geq t'$, for some $t' > 0$, with some $z^{t'} \geq 0$. Here, $\{r_1^t\}$ and $\{r_2^t\}$ are deterministic
613 sequences with $\frac{a_1}{t+1} \leq r_1^t \leq 1$ and $r_2^t \leq \frac{a_2}{(t+1)^\delta}$, with $a_1, a_2 > 0$, and $\delta > 0$. Then, the*

614 following holds: (1) $z^t = O(\frac{1}{t^{\delta-1}})$ provided that $a_1 > \delta - 1$; (2) if $a_1 \leq \delta - 1$, then
 615 $z^t = O(\frac{1}{t^s})$, for any $s < a_1$.

616 We are finally ready to finalize the proof of Theorem 4.6.

617 *Proof.* (Proof of Theorem 4.6) From equations (4.21) and (4.16) we get that
 618 $(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{r}(\mathbf{x}^t)$

$$619 \leq -\min \left\{ \frac{\varphi'_o(0)G_o}{2\|\mathbf{H}\|g'_t} \left(\frac{\lambda_{\mathbf{H}}}{N} - \frac{2S_{\mathbf{H}}}{\sqrt{N}}k \right), \frac{b\varphi'_c(0)G_c}{4agt} \frac{\lambda_2(\mathbf{L})}{1 + \frac{1}{k^2}} \right\} \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2.$$

620 Therefore, taking the expectation in (4.15), we have that

$$622 \mathbb{E}[V(\mathbf{x}^{t+1})] \leq \left(1 - \frac{a_1}{t+1}\right) \mathbb{E}[V(\mathbf{x}^t)] + \frac{a_2}{(1+t)^{2\delta}},$$

623 where

$$625 a_1 = \min \left\{ \frac{\varphi'_o(0)G_o a(1-\delta) (\lambda_{\mathbf{H}} - 2S_{\mathbf{H}}\sqrt{N}k)}{\|\mathbf{H}\|N (\|\mathbf{x}^0 - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\| + b\sqrt{MN}dc_c + a\|\mathbf{H}\|\sqrt{N}c_o)}, \right. \\ \left. \frac{b\varphi'_c(0)G_c(1-\delta)\lambda_2(\mathbf{L})k^2}{2(k^2+1) (\|\mathbf{x}^0\| + b\sqrt{MN}dc_c + a\|\mathbf{H}\|\sqrt{N}c_o)} \right\}$$

626 and $a_2 = a^2c'$. Therefore, using the Lemma 4.11, $\hat{\delta}$ is any positive number such that

$$629 \hat{\delta} < \min \left\{ 2\delta - 1, \frac{\varphi'_o(0)G_o a(1-\delta) (\lambda_{\mathbf{H}} - 2S_{\mathbf{H}}\sqrt{N}k)}{\|\mathbf{H}\|N (\|\mathbf{x}^0 - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\| + b\sqrt{MN}dc_c + a\|\mathbf{H}\|\sqrt{N}c_o)}, \right. \\ \left. \frac{b\varphi'_c(0)G_c(1-\delta)\lambda_2(\mathbf{L})k^2}{2(k^2+1) (\|\mathbf{x}^0\| + b\sqrt{MN}dc_c + a\|\mathbf{H}\|\sqrt{N}c_o)} \right\}.$$

630 Therefore, using Lemma 4.11 we obtain MSE convergence with rate $O(1/t^{\hat{\delta}})$. \square

633 *Remark 4.12.* Even though, we see that the convergence factor $\hat{\delta}$ depends on the
 634 system parameters, i.e., on the network and sensing model and also on the innovation
 635 and consensus nonlinearities, it is easy to see that $\hat{\delta} \in (0, 1)$ regardless of the system
 636 parameters. Recall that Theorem 4.3 shows that the proposed estimator (3.1) obtains
 637 rate $1/t$ in the *weak convergence sense*, while Theorem 4.6 shows that (3.1) obtains
 638 a slower convergence rate, but in the sense of the *mean squared convergence*. Note
 639 that this is not a contradiction, and Theorem 4.6 adds information with respect to
 640 Theorem 4.3. Namely, it is well known that mean squared convergence implies con-
 641 vergence in distribution; therefore, with the same assumptions as in Theorem 4.6, the
 642 convergence rate $1/t^{\hat{\delta}}$ is also attainable for convergence in distribution. In contrast,
 643 from Theorem 4.3, we can not conclude that the rate of the mean squared convergence
 644 is also $1/t$.

645 *Remark 4.13.* In fact, we next show that, in the presence of the heavy-tailed
 646 observation noise considered here, the MSE convergence rate cannot be as fast as $1/t$,
 647 for any estimator (even not for centralized ones). In this sense, the fact that quantity
 648 $\hat{\delta}$ is strictly smaller than one is not a consequence of loose bounds, but it is rather due
 649 to the intrinsic difficulty of the estimation problem. To be specific, we consider here
 650 the special case where each agent i observes a scalar parameter $\theta^* \in \mathbb{R}$ according to
 651 (4.22) $z_i(t) = \theta^* + n_i^t$,
 652 where n_i^t satisfies Assumption 2.2. In this case, the proposed estimator (3.1) can
 653 be viewed as a mean estimator of the probability density function $p_o(u - \theta^*)$. Let us
 654

655 denote by \mathcal{P} the class of all probability density functions $p_o(u - \theta^*)$ such that p_o is the
 656 pdf of the observation noise that satisfies Assumption 2.2, for any $\theta^* \in \mathbb{R}$. Extending
 657 the results from [8] (see Appendix G in [34]), we prove that, for any $\theta^* \in \mathbb{R}$, and for
 658 any mean estimator $\hat{\theta}_t$, the following holds:

$$659 \quad (4.23) \quad \sup_t \sup_{p \in \mathcal{P}} tN \mathbb{E}[|\hat{\theta}_t - \theta^*|^2] = +\infty.$$

660 On the other hand, Theorem 4.6 shows that, with the proposed distributed estima-
 661 tor (3.1), the following holds:

$$663 \quad \sup_t \sup_{p \in \mathcal{P}} (tN)^{\hat{\delta}} \mathbb{E}[|\hat{\theta}_t - \theta^*|^2] < +\infty,$$

664 for some $\hat{\delta} \in (\frac{1}{2}, 1)^3$

666 *Remark 4.14.* Theorems 4.2, 4.3 and 4.6 continue to hold even if the linear trans-
 667 formation vectors \mathbf{h}_i in (2.1) are no longer static (see Appendix H in [34]). That is,
 668 we can allow that each agent i at each time $t = 0, 1, \dots$, makes the observation by:

$$669 \quad (4.24) \quad z_i^t = (\mathbf{h}_i^t)^\top \theta^* + n_i^t.$$

670 Here, for each agent i , for each time step t , the linear transformation vector \mathbf{h}_i^t is a
 671 random variable that satisfies the following assumptions.

- 673 1. For each agent i and each time step $t = 0, 1, \dots$, the linear transformation
 674 vector is given by $\mathbf{h}_i^t = \bar{\mathbf{h}}_i + \tilde{\mathbf{h}}_i^t$, where the vector $\bar{\mathbf{h}}_i \in \mathbf{R}^M$ is deterministic,
 675 and vector $\tilde{\mathbf{h}}_i^t \in \mathbf{R}^M$ is a random vector;
- 676 2. The sequence of vectors $\{\tilde{\mathbf{h}}_1^t, \tilde{\mathbf{h}}_2^t, \dots, \tilde{\mathbf{h}}_N^t\}$ is i.i.d., with finite second moment,
 677 and it is independent of the sequences \mathbf{n}^t and $\boldsymbol{\xi}_{ij}^t$ for $\{i, j\} \in E$;
- 678 3. At each agent $i = 1, \dots, N$ at each time $t = 0, 1, \dots$, each entry $\ell = 1, 2, \dots, M$
 679 $[\tilde{\mathbf{h}}_i^t]_\ell$ has the same probability density function p_h ;
- 680 4. The pdf p_h is symmetric, i.e. $p_h(u) = p_h(-u)$, for every $u \in \mathbb{R}$ and $p_h(u) > 0$
 681 for $|u| \leq c_h$, for some constant $c_h > 0$;
- 682 5. The matrix $\sum_{i=1}^N \bar{\mathbf{h}}_i (\bar{\mathbf{h}}_i)^\top$ is invertible.

683 **5. Analytical and numerical examples.** In this section we provide analytical
 684 and numerical examples that illustrate results from Section 4.

685 **Example 1:** We consider the network where each agent i observes a scalar parameter
 686 $\theta^* \in \mathbb{R}$ following the linear regression model:

$$687 \quad (5.1) \quad z_i(t) = h\theta^* + n_i^t,$$

689 where $h \neq 0$ and $n_i(t)$ is zero mean and i.i.d. in time and across agents. For sim-
 690 plicity, we assume that the underlying graph of the network is regular, with degree
 691 d . We assume that there is no communication noise between agents, i.e., $\boldsymbol{\xi}_{ij} \equiv 0$
 692 for $(i, j) \in E_d$. We additionally assume that the nonlinearity on the consensus part
 693 Ψ_c in (3.1) is the identity function and the nonlinearity on the innovation part is
 694 $\Psi_o(w) = B \tanh(w/B)$, for $B > 0$. Therefore, algorithm (3.1) is now given by:

$$695 \quad (5.2) \quad x_i^{t+1} = x_i^t - \alpha_t \left(\frac{b}{a} \sum_{j \in \Omega_i} (x_i^t - x_j^t) - h\Psi_o(z_i^t - hx_i^t) \right),$$

696 for each agent i and each time t . From Theorem 4.3, we have that the asymptotic
 697 covariance matrix is given by (4.13) and matrix \mathbf{S}_0 is now given by $\mathbf{S}_0 = \sigma_o^2 h^2 \mathbf{I}$
 698 and $\sigma_o^2 = \int |\Psi_o(w)|^2 d\Phi_o(w)$ is the effective observation noise. Following the same
 699 procedure as in [18, 15], for $\boldsymbol{\Sigma} = \frac{1}{2}\mathbf{I} - a\varphi'_o(0)h^2\mathbf{I}$, we have that the average per-agent
 700

³Notice that in the centralized case, the observations are collected in batches of fixed size N . That is, after t time steps, there are Nt observations. Henceforth, we include quantity N in (4.23) for a precise statement. Note that, since N is constant and the supremum is taken with respect to t , the inclusion of N is not necessary.

701 asymptotic variance, denoted by $\sigma_B^2 = \frac{1}{N} \text{Tr}(\mathbf{S})$, is equal to $\sigma_B^2 = \frac{a^2 \sigma_o^2 h^2}{2ah^2 \varphi'_o(0) - 1}$, for
 702 $a > \frac{1}{2h^2 \varphi'_o(0)}$ (see Appendix F in [34]). Therefore, we need to change the constant a
 703 when changing B , i.e., we define $a = a(B) = \frac{1}{2h^2 \varphi'_o(0)(B)} + \epsilon^4$, for some constant $\epsilon > 0$,
 704 we rewrite σ_B^2 as follows (see Appendix F in [34]), $\sigma_B^2 = \frac{(1+2h^2 \varphi'_o(0)\epsilon)^2 \sigma_o^2(B)}{8h^4 \varphi'_o(0)^3 \epsilon}$. For the
 705 nonlinearity Ψ_o that is considered here, we have that $\sigma_o^2 = \int_{-\infty}^{+\infty} B^2 \tanh^2\left(\frac{w}{B}\right) f(w) dw$,
 706 and $\varphi'_o(0) = \int_{-\infty}^{+\infty} \Psi'(w) f(w) dw = \int_{-\infty}^{+\infty} \frac{1}{\cosh^2\left(\frac{w}{B}\right)} f(w) dw$. Notice that both functions
 707 σ_o^2 and $\varphi'_o(0)$ are increasing with respect to B (see Appendix F in [34]). Since we
 708 have that $|B^2 \tanh^2\left(\frac{w}{B}\right) f(w)| \leq |w^2 f(w)|$ and $|\frac{1}{\cosh^2\left(\frac{w}{B}\right)} f(w)| \leq |f(w)|$ for all $w \in \mathbb{R}$
 709 and all $B > 0$, using the Lebesgue's dominated convergence theorem, we have that
 710 $\lim_{B \rightarrow 0^+} \sigma_o^2 = 0$, $\lim_{B \rightarrow +\infty} \sigma_o^2 = \sigma_\eta^2$, $\lim_{B \rightarrow 0^+} \varphi'_o(0) = 0$, $\lim_{B \rightarrow +\infty} \varphi'_o(0) = 1$, where σ_η^2 is the
 711 variance of the observation noise η . Therefore, we have that $\sigma_0^2 = \lim_{B \rightarrow 0^+} \sigma_B^2 = +\infty$
 712 (see Appendix F in [34]), and $\sigma_\infty^2 = \lim_{B \rightarrow +\infty} \sigma_B^2 = \frac{(1+2h^2 \epsilon)^2 \sigma_\eta^2}{8h^4 \epsilon}$. Suppose now that
 713 the variance of the observation noise η is infinite, i.e. $\sigma_\eta^2 = +\infty$. This means that
 714 $\sigma_\infty^2 = +\infty$. For the continuous function σ_B^2 , defined for all $B \in (0, +\infty)$, we have
 715 that $\lim_{B \rightarrow 0^+} \sigma_B^2 = \lim_{B \rightarrow +\infty} \sigma_B^2 = +\infty$. Therefore, there exists an optimal B^* such that
 716 $\sigma_{B^*}^2 = \inf_{B \in (0, \infty)} \sigma_B^2$. Note that the case $B \rightarrow \infty$ corresponds to a \mathcal{LU} scheme from [18],
 717 while the case $B \rightarrow 0$ corresponds to each agent working in isolation. Therefore,
 718 we show analytically on the simple class of nonlinearities Ψ_o (hyperbolic tangent),
 719 that cooperation through a nonlinear mapping Ψ_o strictly improves performance with
 720 respect to both using linear and non-cooperative schemes.

721 To numerically illustrate the above results, we now consider a sensor (agents)
 722 network with $N = 8$ agents, setting that the underlying topology is given by a regular
 723 graph with degree $d = 3$. The true parameter is $\theta^* = 1$, the observation parameter is
 724 $h = 1$, and the observation noise for each agent's measurements has the following pdf

$$\begin{aligned}
 & 725 \quad (5.3) \quad f(w) = \frac{\beta - 1}{2(1 + |w|)^\beta}, \\
 & 726
 \end{aligned}$$

727 with $\beta = 2.05$, which has an infinite variance. Recall that we assumed that there is
 728 no communication noise between agents. We set the consensus parameter as $b = 1$
 729 and the innovation parameter as $a = a(0.3) = \frac{1}{2h^2 \varphi'_o(0)(0.3)} + 0.1$. Figure 1a shows
 730 the average per-agent asymptotic variance σ_B^2 versus B . As it can be seen, optimal
 731 B^* approximately equals $B^* = 0.65$. Using Monte Carlo simulations, we compare
 732 numerically an estimated per-sensor MSE across iterations, for the optimal B^* and
 733 for some sub-optimal choices of B . We can see that the algorithm performs better
 734 for the optimal value B^* than for the other considered suboptimal choices of B (see
 735 Figure 1b), hence confirming the theory.

736 **Example 2:** In this example we provide analysis in the terms of the average
 737 per node variance with respect to the level of the mutual dependence of observation
 738 and communication noise. Once more, we consider the network where each agent i
 739 observes a scalar parameter $\theta^* \in \mathbb{R}$ following the linear regression model (5.1) and we
 740 assume that the underlying graph of the network is regular, with degree d . As it is said,
 741 we now allow observation and communication noise to be mutually dependent. For

⁴ ϵ is added since we need to have that $a > \frac{1}{2h^2 \varphi'_o(0)}$.

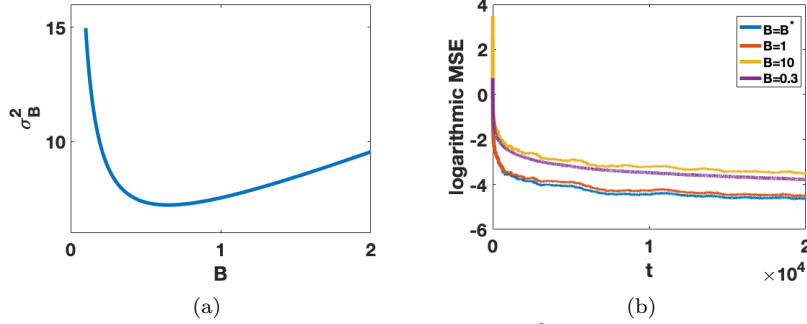


FIG. 1. (a) Average per-agent asymptotic variance σ_B^2 versus B (b) Monte Carlo-estimated per-sensor MSE error on logarithmic scale for the different choices of B

742 simplicity, we consider the case when that dependence between communication noise
 743 ξ_{ij}^t and observation noise n_i is given by $\xi_{ij} = \rho n_i^t + \sqrt{1 - \rho^2} \hat{n}_i^t$, at each time $t = 0, 1, \dots$
 744 and for all tuples $\{i, j\} \in E$, where, $\rho \in (-1, 1)$, sequence $\{\hat{n}_i^t\}$ is independently
 745 identically distributed in time t and across all agents i . Moreover, n_i^t are \hat{n}_j^s mutually
 746 independent whenever $(i, t) \neq (j, s)$. Here, it is easy to see that we have strong positive
 747 correlation if $\rho \rightarrow 1$, strong negative correlation if $\rho \rightarrow -1$ and we do not have any
 748 correlation if $\rho = 0$. Moreover, we set that $\Psi_o(w) = \Psi_c(w) = \text{sign } w$, and hence,
 749 algorithm (3.1) is given by

(5.4)

$$750 \quad x_i^{t+1} = x_i^t - \alpha_t \left(\frac{b}{a} \sum_{j \in \Omega_i} \Psi_c \left(x_i^t - x_j^t + \rho n_i^t + \sqrt{1 - \rho^2} \hat{n}_i^t \right) - h \Psi_o \left(h(\theta^* - x_i^t) + n_i \right) \right).$$

751 Analogously to the previous example, we have that the average per-agent asymptotic
 752 variance σ_ρ^2 is given by

$$754 \quad (5.5) \quad \sigma_\rho^2 = \frac{b^2 \sigma_c^2 d^2 + a^2 h^2 \sigma_o^2 - 2abhd\sigma_{oc}}{N(2ah^2\varphi'_o(0) - 1)}$$

$$755 \quad (5.6) \quad + \frac{b^2 \sigma_c^2 d^2 + a^2 h^2 \sigma_o^2 - 2abhd\sigma_{oc}}{N} \sum_{i=2}^N \frac{1}{2b\varphi'_c(0)\lambda_i + 2ah^2\varphi'_o(0) - 1},$$

756 since $\mathbf{S}_0 = \left(\frac{b^2}{a^2} \sigma_c^2 d^2 + \sigma_o^2 h^2 - 2\frac{b}{a} hd\sigma_{oc} \right) \mathbf{I}$ and $\mathbf{\Sigma} = \frac{1}{2} \mathbf{I} - a \left(\frac{b}{a} \varphi'_c(0) \mathbf{L} + \varphi'_o(0) h^2 \mathbf{I} \right)$.
 757 Here, regardless of ρ we have that $\sigma_o^2 = \sigma_c^2 = 1$ and $\varphi'_o(0) = 2p_n(0)$ (see [15]). On the
 758 other hand, σ_{oc} which is effective cross-covariance between the observation and the
 759 communication noise after passing through the appropriate nonlinearity and $\varphi'_c(0)$
 760

761 are functions with respect to ρ . We have that

$$762 \quad (5.7) \quad \sigma_{\text{oc}} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \Psi_c(\rho x + \sqrt{1-\rho^2}y) \Psi_o(x) p_{\hat{n}}(y) p_n(x) dx dy$$

$$763 \quad (5.8) \quad = \int_0^{+\infty} \int_{\frac{-\rho x}{\sqrt{1-\rho^2}}}^{\infty} p_{\hat{n}}(y) p_n(y) dy dx - \int_0^{+\infty} \int_{-\infty}^{\frac{-\rho x}{\sqrt{1-\rho^2}}} p_{\hat{n}}(y) p_n(y) dy dx$$

$$764 \quad (5.9) \quad - \int_{-\infty}^0 \int_{\frac{-\rho x}{\sqrt{1-\rho^2}}}^{\infty} p_{\hat{n}}(y) p_n(y) dy dx + \int_{-\infty}^0 \int_{-\infty}^{\frac{-\rho x}{\sqrt{1-\rho^2}}} p_{\hat{n}}(y) p_n(y) dy dx,$$

765 and we see that $\sigma_{\text{oc}} \rightarrow 0$ as $\rho \rightarrow 0$, $\sigma_{\text{oc}} \rightarrow 1$ as $\rho \rightarrow 1$ and $\sigma_{\text{oc}} \rightarrow -1$ as $\rho \rightarrow -1$.

767 Moreover, we have that $\varphi'_c(0) = 2 \int_{-\infty}^{\infty} p_{\hat{n}}(-\rho x) p_n(\sqrt{1-\rho^2}x) dx$, and again, it is easy

768 to see that, $\varphi'_c(0) \rightarrow 2p_n(0)$ as $\rho \rightarrow \pm 1$ and $\varphi'_c(0) \rightarrow 2p_{\hat{n}}(0)$ as $\rho \rightarrow 0$. To demonstrate

769 the above results, again we consider a sensor (agents) network with $N = 8$ agents,

770 setting that the underlying topology is given by a regular graph with degree $d = 3$.

771 The true parameter is $\theta^* = 1$, the observation parameter, the innovation parameter

772 and consensus parameter are $h = a = b = 1$. We set that for all i , n_i and \hat{n}_i have

773 the pdf as in (5.3) with $\beta = 2.05$. Figure 2a shows σ_{ρ}^2 with respect to ρ . As it can be

774 seen, the lowest σ_{ρ}^2 is attained at $\rho = 1$, also σ_{ρ}^2 has two local maxima at $\rho \approx -0.88$

775 and at $\rho \approx 0.31$. Figure 2b shows the comparison of Monte Carlo simulation for

776 $\frac{1}{N} \|\mathbf{x}^t - \mathbf{1} \otimes \theta\|^2 t$ for different choices of ρ . Moreover, Figure 2b justifies the results

777 presented in 2a, in the sense that $\frac{1}{N} \|\mathbf{x}^t - \mathbf{1} \otimes \theta\|^2 t$ is minimal for $\rho = 1$ and maximal

778 for $\rho = -0.88$. Finally, we note that, while the two local maxima obtained here are

779 specific for the simplistic correlation and sensing model assumed here for analytical

780 tractability, we observe numerically for more general models that the general trend of

781 this example is preserved, in the sense that higher (more positive) correlations lead

782 to a better performance.

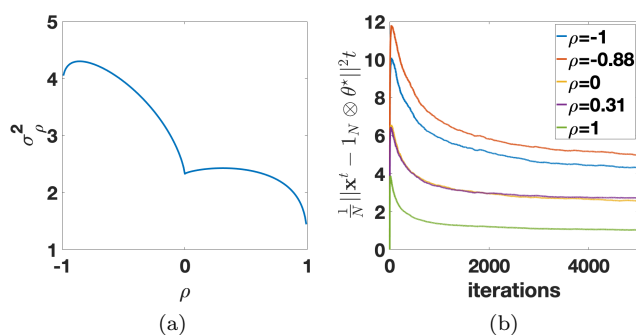


FIG. 2. (a) Average per-agent asymptotic variance σ_B^2 versus B (b) Monte Carlo-estimation of $\frac{1}{N} \|\mathbf{x}^t - \mathbf{1}_N \otimes \theta^*\|^2 t$ for different choices of ρ

783 **5.1. Numerical simulations.** In this subsection, we demonstrate the perfor-
 784 mance of proposed consensus+innovations estimator in a larger sensor network. We
 785 consider a sensor network with $N = 40$ agents where the underlying topology is an
 786 instance of a random geometric graph; we used randomly generated true parameter
 787 $\theta^* \in \mathbb{R}^{10}$, whose entries are drawn mutually independently from the uniform distri-

788 bution on $[-10, 10]$; we used randomly generated observation vectors $\mathbf{h}_i \in \mathbb{R}^{10}$, for
 789 which the condition 2 of Assumption 2.1 is verified to be true. We set the consensus
 790 parameter as $b = 1$ and step-size parameter as $\delta = 1$. First, we compare the
 791 proposed consensus+innovations estimator with the method from [1] and its hypo-
 792 theoretical variant in the case when there is no communication noise, but in the presence
 793 of heavy-tailed observation noise with pdf as in (5.3) for $\beta = 2.05$. Here, we used
 794 the same algorithm settings and the same nonlinearities for the proposed algorithm
 795 as in Example 1, with a slight change, i.e., we set that $B = 10$ and $a = 0.2$. For
 796 method from [1] and its hypothetical variant (see Appendix E in [34]), we set that
 797 $B_i = 2$, $\phi_{i,1}(x) = x$ and $\phi_{i,2}(x) = \tanh(x)$ for all agents i . Furthermore, we set
 798 that weighting coefficients are chosen according to $a_{ij} = \frac{\tilde{\mathbf{A}}_{ij}}{\sum_{\ell \in \mathcal{N}_i} \mathbf{A}_{\ell i}}$, where $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$.

799 Moreover, for the smoothing recursions, zero initial conditions are assumed, ν_i is set
 800 to 0.9 for every agent i and $\epsilon = 10^{-2}$. We can see all methods manage to (slowly)
 801 decrease MSE over iterations, with the proposed method exhibiting the best perfor-
 802 mance among the three methods considered. Figure 3b shows Monte Carlo simulation
 803 of the MSE for the proposed algorithm, algorithm from [1] and the algorithm in [15],
 804 when communication between agents is also contaminated with heavy-tailed commu-
 805 nication noise. Here, for the proposed algorithm we set that both nonlinearities are
 806 $\Psi_o(w) = \Psi_c(w) = B \tanh(w/B)$, for $B = 10$ and $a = 1$. Further, we use the same
 807 algorithm setting for the method in [1] as in the previous simulation example, and
 808 we use the same nonlinearity on the consensus part and the same B for algorithm
 809 from [15] as in the proposed algorithm. We can see that both [15] and [1] here fail to
 810 converge, while the proposed method still effectively reduces MSE.

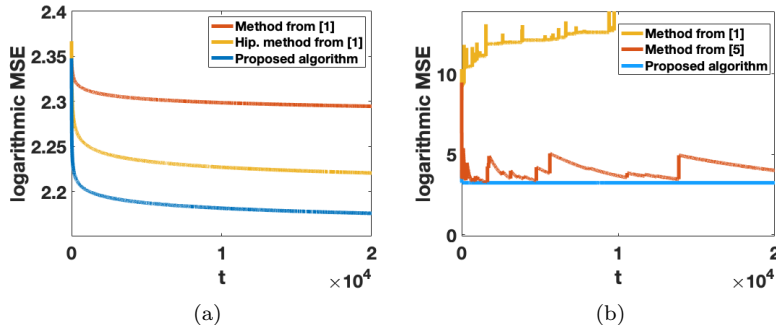


FIG. 3. (a) Monte Carlo-estimated per-sensor MSE error on logarithmic scale for proposed algorithm for $B = 10$, method from [1] and its hypothetical variant (b) Monte Carlo-estimated per-sensor MSE error on logarithmic scale for proposed algorithm, algorithm form [1] and algorithm from [15]

811 We next present the scenario where the observation and communication noises are
 812 mutually dependent. To do this, we set that the i -th element of the observation noise
 813 \mathbf{n} is given by $\mathbf{n}_i = \mathbf{v}_i \exp(\frac{h}{2}\mathbf{v}_i^2)$, where \mathbf{v} has standard normal distribution and h is a
 814 heavy-tail parameter (see [9]). Moreover, the ℓ -th element of the communication noise
 815 ξ_{ij} is given by $[\xi_{ij}]_\ell = [\mathbf{w}_{ij}]_\ell \exp(\frac{h}{2}[\mathbf{w}_{ij}]_\ell^2)$, where \mathbf{w}_{ij} is the linear transformation of
 816 \mathbf{v} , i.e., $\mathbf{w}_{ij} = \mathbf{W}_{ij}\mathbf{v}$ and $\mathbf{W}_{ij} \in \mathbb{R}^{M \times N}$ is a randomly generated matrix independent
 817 of the observation noise. Figure 4a presents Monte Carlo estimates of per-agent
 818 MSE across iterations. Figure 4b shows Monte Carlo simulation of quantity $\frac{1}{N}\|\mathbf{x}^t -$
 819 $\mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2 \sqrt{t}$. For this numerical setting, from the Figure 4b, we can deduce that
 820 $E[\|\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2]$ decreases at least as fast as $O(\frac{1}{\sqrt{t}})$, hence confirming our MSE
 821 rate theory.

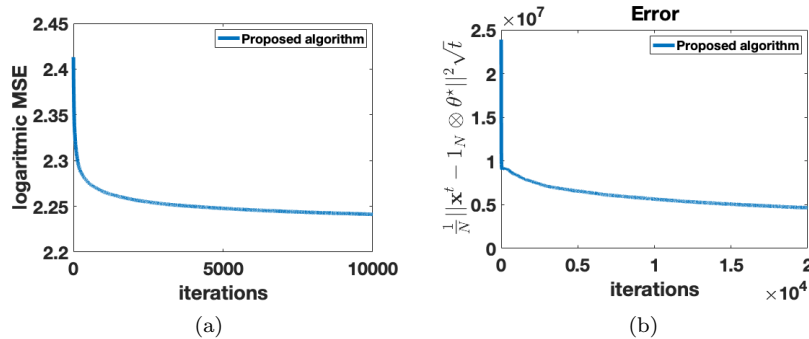


FIG. 4. (a) Monte Carlo-estimated per-sensor MSE error on logarithmic scale for proposed algorithm when link failures can occur for $B = 1$ and $h = 10$ (b) Monte Carlo-estimation of $\frac{1}{N} \|\mathbf{x}^t - \mathbf{1}_N \otimes \theta^*\|^2 \sqrt{t}$ for $B = 10$ and $h = 2$.

822

6. Conclusion.

823

We have studied distributed consensus+innovations estimation under the simultaneous presence of heavy-tailed (infinite variance) correlated sensing and communication noises. This setting is in contrast with existing work that either always assumes a finite-variance sensing noise. We developed a nonlinear estimator and established its almost sure convergence and asymptotic normality. Furthermore, we showed that the estimator achieves a sublinear MSE convergence rate $O(1/t^\kappa)$, and we explicitly characterized the rate $\kappa \in (0, 1)$ in terms of system parameters. Analytical examples illustrate the role of the nonlinearities incorporated in the method and the effects of noises correlation. Finally, numerical simulations corroborate our findings and demonstrate that the proposed distributed estimator converges under the simultaneous presence of heavy-tailed (infinite variance) correlated sensing and communication noises, while, for the same setting, existing distributed estimators fail to converge.

835

Acknowledgements

836

The work of D. Bajovic, D. Jakovetic and M. Vukovic is supported by the Ministry of Education, Science and Technological Development, Republic of Serbia. Moreover, the work of D. Bajovic and D. Jakovetic is supported by the European Union's Horizon 2020 Research and Innovation program under grant agreement No. 957337. The paper reflects only the view of the authors and the Commission is not responsible for any use that may be made of the information it contains. The work of D. Jakovetic is also supported by the Provincial Secretariat for Higher Education and Scientific Research, grant no 142-451-2593/2021-01/2.

843

REFERENCES

844

- [1] S. AL-SAYED, A. M. ZOUBIR, AND A. H. SAYED, *Robust distributed estimation by networked agents*, IEEE Transactions on Signal Processing, 65 (2017), pp. 3909–3921.
- [2] D. BAJOVIC, D. JAKOVETIC, J. M. MOURA, J. XAVIER, AND B. SINOPOLI, *Large deviations performance of consensus+ innovations distributed detection with non-gaussian observations*, IEEE Transactions on Signal Processing, 60 (2012), pp. 5987–6002.
- [3] D. BAJOVIC, D. JAKOVETIC, J. XAVIER, B. SINOPOLI, AND J. MOURA, *Distributed detection via gaussian running consensus: Large deviations asymptotic analysis*, Signal Processing, IEEE Transactions on, 59 (2011), pp. 4381 – 4396.
- [4] Y. CHEN, S. KAR, AND J. MOURA, *Resilient distributed field estimation*, SIAM Journal on Control and Optimization, 58 (2020), pp. 1429–1456.
- [5] S. CHOUVARDAS, K. SLAVAKIS, AND S. THEODORIDIS, *Adaptive robust distributed learning in diffusion sensor networks*, Signal Processing, IEEE Transactions on, 59 (2011), pp. 4692 – 4707.
- [6] L. CLAVIER, T. PEDERSEN, I. LARRAD, M. LAURIDSEN, AND M. EGAN, *Experimental evidence*

857

- 858 for heavy tailed interference in the IoT, IEEE Communications Letters, 25 (2021), pp. 692–
859 695.
- 860 [7] S. DASARATHAN, C. TEPEDELENLIOGLU, M. BANAVAR, AND A. SPANIAS, *Robust consensus in the*
861 *presence of impulsive channel noise*, Signal Processing, IEEE Transactions on, 63 (2014).
- 862 [8] L. DEVROYE, M. LERASLE, G. LUGOSI, AND R. I. OLIVEIRA, *Sub-gaussian mean estimators*,
863 The Annals of Statistics, 44 (2016), pp. 2695–2725.
- 864 [9] G. GOERG, *The lambert way to gaussianize heavy-tailed data with the inverse of tukey’s h*
865 *transformation as a special case*, The Scientific World Journal, 2015 (2015).
- 866 [10] M. HAENGGI AND R. GANTI, *Interference in large wireless networks*, Foundations and Trends
867 in Networking, 3 (2009), pp. 127–248.
- 868 [11] B. HUGHES, *Alpha-stable models of multiuser interference*, in 2000 IEEE International Sympo-
869 sium on Information Theory (Cat. No.00CH37060), 2000, pp. 383–.
- 870 [12] J. ILOW AND D. HATZINAKOS, *Analytic alpha-stable noise modeling in a poisson field of in-*
871 *terferers or scatterers*, Signal Processing, IEEE Transactions on, 46 (1998), pp. 1601 –
872 1611.
- 873 [13] D. JAKOVETIC, D. BAJOVIC, A. K. SAHU, S. KAR, N. MILOSEVIC, AND D. STAMENKOVIC, *Nonlin-*
874 *ear gradient mappings and stochastic optimization: A general framework with applications*
875 *to heavy-tail noise*. arXiv preprint, <https://arxiv.org/abs/2204.02593>, 2022.
- 876 [14] D. JAKOVETIC, J. M. F. MOURA, AND J. XAVIER, *Distributed detection over noisy networks:*
877 *Large deviations analysis*, IEEE Transactions on Signal Processing, 60 (2012), pp. 4306–
878 4320.
- 879 [15] D. JAKOVETIC, M. VUKOVIC, D. BAJOVIC, A. K. SAHU, AND S. KAR, *Distributed recursive*
880 *estimation under heavy-tail communication noise*, SIAM Journal on Control and Opti-
881 mization, 61 (2023), pp. 1582–1609.
- 882 [16] S. KAR AND J. MOURA, *Asymptotically efficient distributed estimation with exponential family*
883 *statistics*, IEEE Transactions on Information Theory, 60 (2014), pp. 4811–4831.
- 884 [17] S. KAR, J. MOURA, AND H. V. POOR, *Distributed linear parameter estimation: Asymptoti-*
885 *cally efficient adaptive strategies*, SIAM Journal on Control and Optimization, 51 (2013),
886 pp. 2200–2229.
- 887 [18] S. KAR, J. M. F. MOURA, AND K. RAMANAN, *Distributed parameter estimation in sensor*
888 *networks: Nonlinear observation models and imperfect communication*, IEEE Transactions
889 on Information Theory, 58 (2012), pp. 3575–3605.
- 890 [19] S. KUMAR, U. K. SAHOO, A. K. SAHOO, AND D. P. ACHARYA, *Diffusion minimum-wilcoxon-*
891 *norm over distributed adaptive networks: Formulation and performance analysis*, Digital
892 Signal Processing, 51 (2016), pp. 156–169.
- 893 [20] A. LALITHA, T. JAVIDI, AND A. D. SARWATE, *Social learning and distributed hypothesis testing*,
894 IEEE Transactions on Information Theory, 64 (2018), pp. 6161–6179.
- 895 [21] Z. LI AND S. GUAN, *Diffusion normalized huber adaptive filtering algorithm*, Journal of the
896 Franklin Institute, 355 (2018), pp. 3812–3825.
- 897 [22] Q. LIU AND A. IHLER, *Distributed estimation, information loss and exponential families*. arXiv
898 preprint, <https://arxiv.org/abs/1410.2653>, 2014.
- 899 [23] C. LOPES AND A. SAYED, *Diffusion least-mean squares over adaptive networks: Formulation*
900 *and performance analysis*, IEEE Transactions on Signal Processing, 56 (2008), pp. 3122–
901 3136.
- 902 [24] G. MATEOS, I. SCHIZAS, AND G. GIANNAKIS, *Distributed recursive least-squares for consensus-*
903 *based in-network adaptive estimation*, Signal Processing, IEEE Transactions on, 57 (2009),
904 pp. 4583 – 4588.
- 905 [25] V. MATTA, P. BRACA, S. MARANO, AND A. H. SAYED, *Diffusion-based adaptive distributed*
906 *detection: Steady-state performance in the slow adaptation regime*, IEEE Transactions on
907 Information Theory, 62 (2016), pp. 4710–4732.
- 908 [26] S. MODALAVALASA, U. SAHOO, A. SAHOO, AND S. BARAHA, *A review of robust distributed esti-*
909 *mation strategies over wireless sensor networks*, Signal Processing, 188 (2021), p. 108150.
- 910 [27] A. NEDIC, A. OLSHEVSKY, AND C. A. URIBE, *Nonasymptotic convergence rates for cooperative*
911 *learning over time-varying directed graphs*, in 2015 American Control Conference (ACC),
912 IEEE, 2015, pp. 5884–5889.
- 913 [28] M. B. NEVEL’SON AND R. Z. HAS’ MINSKII, *Stochastic approximation and recursive estimation*,
914 vol. 47, American Mathematical Soc., 1976.
- 915 [29] B. POLYAK AND Y. TSYPKIN, *Adaptive estimation algorithms: Convergence, optimality, stabil-*
916 *ity*, Automation and Remote Control, 1979 (1979).
- 917 [30] S. RAM, V. VEERAVALLI, AND A. NEDIC, *Distributed and Recursive Parameter Estimation*,
918 Springer Science & Business Media, 2009, pp. 17–38.
- 919 [31] B. SELIM, M. S. ALAM, V. CARVALHO, G. KADDOUM, AND B. L. AGBA, *Noma-based iot net-*

- 920 *works: Impulsive noise effects and mitigation*, IEEE Communications Magazine, 58 (2020),
921 pp. 69–75.
- 922 [32] W. SHI, Q. LING, G. WU, AND W. YIN, *Extra: An exact first-order algorithm for decentralized*
923 *consensus optimization*. arXiv preprint, <https://arxiv.org/abs/1404.6264>, 2014.
- 924 [33] S. THEODORIDIS, K. SLAVAKIS, AND I. YAMADA, *Adaptive learning in a world of projections*,
925 Signal Processing Magazine, IEEE, 28 (2011), pp. 97 – 123.
- 926 [34] M. VUKOVIC, D. JAKOVETIC, D. BAJOVIC, AND S. KAR, *Nonlinear consensus+innovations*
927 *under correlated heavy-tailed noises: Mean square convergence rate and asymptotics*. arXiv
928 preprint, <https://arxiv.org/abs/2212.11959>, 2022.
- 929 [35] F. WEN, *Diffusion least mean p -power algorithms for distributed estimation in alpha-stable*
930 *noise environments*, Electronics Letters, 49 (2013).
- 931 [36] M. WIN, P. PINTO, AND L. SHEPP, *A mathematical theory of network interference and its*
932 *applications*, Proceedings of the IEEE, 97 (2009), pp. 205 – 230.
- 933 [37] X. YANG AND A. PETROPULU, *Co-channel interference modeling and analysis in a poisson*
934 *field of interferers in wireless communications*, IEEE Transactions on Signal Processing,
935 51 (2003), pp. 64–76.
- 936 [38] X. ZHAO, S.-Y. TU, AND A. H. SAYED, *Diffusion adaptation over networks under imperfect*
937 *information exchange and non-stationary data*, IEEE Transactions on Signal Processing,
938 60 (2012), pp. 3460–3475.

939 **Appendix.**

940 **A. Some results on Stochastic approximation.** We make use of the following
 941 standard stochastic approximation result, see [28], see also [18].

942 **THEOREM 6.1.** Let $\{\mathbf{x}^t \in \mathbb{R}^l\}_{t \geq 0}$ be a random sequence:

$$943 \quad (6.1) \quad \mathbf{x}^{t+1} = \mathbf{x}^t + \alpha_t [\mathbf{r}(\mathbf{x}^t) + \boldsymbol{\gamma}(t+1, \mathbf{x}^t, \omega)],$$

945 where, $\mathbf{r}(\cdot) : \mathbb{R}^l \rightarrow \mathbb{R}^l$ is Borel measurable and $\{\boldsymbol{\gamma}(t, \mathbf{x}, \omega)\}_{t \geq 0, \mathbf{x} \in \mathbb{R}^l}$ is a family of
 946 random vectors in \mathbb{R}^l , defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, and $\omega \in \Omega$ is a canonical
 947 element. Let the following sets of assumptions hold:

948 **B1:** The function $\boldsymbol{\gamma}(t, \cdot, \cdot) : \mathbb{R}^l \times \Omega \rightarrow \mathbb{R}^l$ is $\mathcal{B}^l \otimes \mathcal{F}$ measurable for every t ; \mathcal{B}^l is
 949 the Borel algebra of \mathbb{R}^l .

950 **B2:** There exists a filtration $\{\mathcal{F}_t\}_{t \geq 0}$ of \mathcal{F} , such that, for each t , the family of
 951 random vectors $\{\boldsymbol{\gamma}(t, \mathbf{x}, \omega)\}_{\mathbf{x} \in \mathbb{R}^l}$ is \mathcal{F}_t measurable, zero-mean and independent
 952 of \mathcal{F}_{t-1} .

953 (If Assumptions B1, B2 hold, $\{\mathbf{x}(t)\}_{t \geq 0}$, is Markov.)

954 **B3:** There exists a twice continuously differentiable $V(\mathbf{x})$ with bounded second
 955 order partial derivatives and a point $\mathbf{x}^* \in \mathbb{R}^l$ satisfying

$$956 \quad V(\mathbf{x}^*) = 0, V(\mathbf{x}) > 0, \mathbf{x} \neq \mathbf{x}^*, \lim_{\|\mathbf{x}\| \rightarrow \infty} V(\mathbf{x}) = \infty,$$

$$957 \quad \sup_{\epsilon < \|\mathbf{x} - \mathbf{x}^*\| < \frac{1}{\epsilon}} \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle < 0, \forall \epsilon > 0.$$

958 **B4:** There exists constants $k_1, k_2 > 0$, such that,

$$959 \quad \|\mathbf{r}(\mathbf{x})\|^2 + \mathbb{E}[\|\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)\|^2] \leq k_1(1 + V(\mathbf{x})) - k_2 \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle$$

960 **B5:** The weight sequence $\{\alpha(t)\}_{t \geq 0}$ satisfies

$$961 \quad \alpha_t > 0, \sum_{t \geq 0} \alpha_t = \infty, \sum_{t \geq 0} \alpha_t^2 < \infty.$$

962 **C1:** The function $\mathbf{r}(\mathbf{x})$ admits the representation

$$963 \quad (6.2) \quad \mathbf{r}(\mathbf{x}) = \mathbf{B}(\mathbf{x} - \mathbf{x}^*) + \boldsymbol{\delta}(\mathbf{x}),$$

964 where

$$965 \quad (6.3) \quad \lim_{\mathbf{x} \rightarrow \mathbf{x}^*} \frac{\|\boldsymbol{\delta}(\mathbf{x})\|}{\|\mathbf{x} - \mathbf{x}^*\|} = 0.$$

966 (Note, in particular, if $\boldsymbol{\delta}(\mathbf{x}) \equiv 0$ then (6.3) is satisfied.)

967 **C2:** The weight sequence $\{\alpha_t\}_{t \geq 0}$ is of form

$$968 \quad (6.4) \quad \alpha_t = \frac{a}{t+1}, \forall t \geq 0,$$

969 where $a > 0$ is a constant (note that **C2** implies **B5**).

970 **C3:** Let \mathbf{I} be the $l \times l$ identity matrix and a, \mathbf{B} as in (6.4) and (6.2), respectively.

971 Then, the matrix $\boldsymbol{\Sigma} = a\mathbf{B} + \frac{1}{2}\mathbf{I}$ is stable.

972 **C4:** The entries of the matrices, $\forall t \geq 0, x \in \mathbb{R}^l$,

$$973 \quad \mathbf{A}(t, \mathbf{x}) = \mathbb{E}[\boldsymbol{\gamma}(t, \mathbf{x}, \omega) \boldsymbol{\gamma}^\top(t, \mathbf{x}, \omega)],$$

974 are finite, and the following limit exists:

$$975 \quad \lim_{t \rightarrow \infty, \mathbf{x} \rightarrow \mathbf{x}^*} \mathbf{A}(t, \mathbf{x}) = \mathbf{S}_0.$$

976 **C5:** There exists $\epsilon > 0$, such that

$$977 \quad \lim_{R \rightarrow \infty} \sup_{\|\mathbf{x} - \mathbf{x}^*\| < \epsilon} \sup_{t \geq 0} \int_{\|\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)\| > R} \|\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)\|^2 dP = 0$$

978 Let Assumptions B1–B5 hold for $\{\mathbf{x}(t)\}_{t \geq 0}$ in (6.1). Then, starting from an arbitrary
 979 initial state, the Markov process, $\{\mathbf{x}^t\}_{t \geq 0}$, converges a.s. to \mathbf{x}^* . In other words,

$$980 \quad \mathbf{P}[\lim_{t \rightarrow \infty} \mathbf{x}^t = \mathbf{x}^*] = 1.$$

981 The normalized process, $\{\sqrt{t}(\mathbf{x}^t - \mathbf{x}^*)\}_{t \geq 0}$, is asymptotically normal if, besides As-

992 *sumptions B1–B5, Assumptions C1–C5 are also satisfied. In particular, as $t \rightarrow \infty$*
 993 (6.5) $\sqrt{t}(\mathbf{x}^t - \mathbf{x}^*) \Rightarrow \mathcal{N}(\mathbf{0}, \mathbf{S}),$

995 *where \Rightarrow denotes convergence in distribution (weak convergence). Also, asymptotic*
 996 *variance, \mathbf{S} , in (6.5) is*

$$997 \quad \mathbf{S} = a^2 \int_0^\infty e^{\Sigma v} \mathbf{S}_0 e^{\Sigma^\top v} dv$$

999 **B. Additional results on nonlinearity φ .** We present some properties of the
 1000 function φ defined in (4.1). As it is stated in [15], we can intuitively see function φ as
 1001 a convolution-like transformation of nonlinearity $\Psi : \mathbb{R} \rightarrow \mathbb{R}$, where the convolution
 1002 is taken with respect to the probability density function p of random value w . If w is
 1003 generated by the underlying probability space $(\Omega, \mathcal{F}, \mathbb{P})$, we have that expectation of
 1004 (6.6) $v = \Psi(a + w) - \varphi(a)$

1005 is equal to zero, i.e., $\mathbb{E}[v] = 0$. Here, the expectation is taken with respect to \mathcal{F} .
 1006 Hence, for all $t = 0, 1, \dots$, we have that expectation of both of the sequences ζ^t, η^t
 1007 defined in (4.3) is equal to zero, due to the fact that communication noise ξ^t and
 1008 observation noise $\mathbf{n}^t, t = 0, 1, \dots$, are generated by underlying probability space.
 1009 We have following Lemma (see [29], see also [15]).

1011 **LEMMA 6.2 ([29]).** *Consider function φ in (4.1), where function $\Psi : \mathbb{R} \rightarrow \mathbb{R}$,*
 1012 *satisfies Assumption 3.1. Then, the following holds:*

- 1013 1. φ is odd;
- 1014 2. If $|\Psi(v)| \leq c_1$, for any $v \in \mathbb{R}$, then $|\varphi(a)| \leq c'_1$, for any $a \in \mathbb{R}$, for some
- 1015 $c'_1 > 0$;
- 1016 3. $\varphi(a)$ is monotonically nondecreasing;
- 1017 4. $\varphi(a) > 0$, for any $a > 0$.
- 1018 5. φ is continuous at zero;
- 1019 6. φ is differentiable at zero, with a strictly positive derivative at zero, equal to:

$$1020 \quad (6.7) \quad \varphi'(0) = \sum_{i=1}^s (\Psi(\nu_i + 0) - \Psi(\nu_i - 0)) p(\nu_i) + \sum_{i=0}^s \int_{\nu_i}^{\nu_{i+1}} \Psi'(v) p(v) dv,$$

1021 *where $\nu_i, i = 1, \dots, s$ are points of discontinuity of Ψ such that $\nu_0 = -\infty$ and*
 1022 *$\nu_{s+1} = +\infty$, and we recall that $p(u)$ is the pdf of random variable w .*

1023 From Lemma 6.2, we have that $\varphi(a) = 0$ if and only if $a = 0$. Moreover, there exists
 1024 a function $\delta : \mathbb{R} \rightarrow \mathbb{R}$, which is continuous in the vicinity of zero, such that

$$1025 \quad (6.8) \quad \varphi(a) = \varphi(0) + \varphi'(0)a + \delta(a) = \varphi'(0)a + \delta(a),$$

1027 and $\lim_{a \rightarrow 0} \frac{\delta(a)}{a} = 0$.

1028 We now prove boundedness of the function $\mathbf{r}(\cdot)$ in equation (4.10). If condition 2
 1029 of Lemma 6.2 is satisfied for both functions φ_c and φ_o , then the right hand side
 1030 of (4.10) would be lesser or equal to some positive constant c , which would led to
 1031 $\|\mathbf{r}(\mathbf{x})\|^2 \leq c_1(1 + V(\mathbf{x}))$. Suppose now that condition 3 of Lemma 6.2 is satisfied for

1032 the function φ_c , then there exists some positive constant c_1 such that

$$\begin{aligned}
1033 \quad \left\| \frac{b}{a} \mathbf{L}_{\varphi_c}(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) \right\|^2 &= \left(\frac{b}{a} \right)^2 \sum_{i=1}^N \left\| \sum_{j \in \Omega_i} \varphi_c(\mathbf{x}_i - \mathbf{x}_j) \right\|^2 \\
1034 &\leq \left(\frac{b}{a} \right)^2 \sum_{i=1}^N \sum_{j \in \Omega_i} \|\varphi_c(\mathbf{x}_i - \mathbf{x}_j)\|^2 \\
1035 &\leq \left(\frac{b}{a} \right)^2 \sum_{i=1}^N \sum_{j \in \Omega_i} \left(c \left(1 + \|\mathbf{x}_i - \mathbf{x}_j\|^2 \right) \right) \\
1036 &\leq \left(\frac{b}{a} \right)^2 \sum_{i=1}^N \sum_{j \in \Omega_i} \left(c \left(1 + \|\mathbf{x}_i - \boldsymbol{\theta}^*\|^2 + \|\mathbf{x}_j - \boldsymbol{\theta}^*\|^2 \right) \right) \\
1037 &\leq c_1 (1 + V(\mathbf{x})),
\end{aligned}$$

1038 since we have that $\|\mathbf{x}_i - \boldsymbol{\theta}^*\|^2 \leq \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2 = V(\mathbf{x})$ for all $i = 1, 2, \dots, N$. If we
1039 assume that condition 3 of Lemma 6.2 is satisfied for the function φ_o , we will get that

$$\begin{aligned}
1040 \quad \left\| \mathbf{H}^\top \varphi_o(\mathbf{H}(\mathbf{x} - (\mathbf{1}_N \otimes \boldsymbol{\theta}^*))) \right\|^2 &\leq \|\mathbf{H}\|^2 \|\varphi_o(\mathbf{H}(\mathbf{x} - (\mathbf{1}_N \otimes \boldsymbol{\theta}^*)))\|^2 \\
1041 &\leq \|\mathbf{H}\|^2 c \left(1 + \|\mathbf{H}(\mathbf{x} - (\mathbf{1}_N \otimes \boldsymbol{\theta}^*))\|^2 \right) \\
1042 &\leq \|\mathbf{H}\|^2 c \left(1 + \|\mathbf{H}\|^2 \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2 \right).
\end{aligned}$$

1043 Therefore, $\left\| \mathbf{H}^\top \varphi_o(\mathbf{H}(\mathbf{x} - (\mathbf{1}_N \otimes \boldsymbol{\theta}^*))) \right\|^2 \leq c_1 (1 + V(\mathbf{x}))$, for some positive constant
1044 c_1 . Hence, inequality in (4.10) is proven.

1045 Next we prove boundedness of $\mathbb{E} \left[\|\gamma(t+1, \mathbf{x}^t, \omega)\|^2 \right]$ in (4.12). If the function Ψ
1046 in (4.1) satisfies condition 5' of Assumption 3.1, whether w in (4.1) has finite or
1047 infinite variance, v in (6.6) is bounded, i.e.,

$$1048 \quad |v|^2 \leq |\Psi(a+w)|^2 + |\varphi(a)|^2 \leq c,$$

1049 for some positive constant c . If the function Ψ in (4.1) satisfies condition 5 of Assumption
1050 3.1 and w has finite variance, we get that variance of v in (6.6) is bounded
1051 with $c(1 + |a|^2)$ for some positive constant c , i.e.,

$$\begin{aligned}
1052 \quad \mathbb{E}[|v|^2] &\leq \mathbb{E}[|\Psi(a+w)|^2 + |\varphi(a)|^2] \leq \mathbb{E}[c_1(1 + |a+w|^2) + c'_1(1 + |a|^2)] \\
1053 &\leq c_1(1 + |a|^2 + \mathbb{E}[|w|^2]) + c'_1(1 + |a|^2) \leq c(1 + |a|^2),
\end{aligned}$$

1054 where c_1 and c_2 are some positive constants. Thus, whether condition 5 or 5' is
1055 satisfied for the function Ψ in (4.1), variance of v in (6.6) is bounded with $c(1 + |a|^2)$
1056 for some positive constant c . Hence, we have that for $\boldsymbol{\zeta}^t, \boldsymbol{\eta}^t$ defined in (4.3)

$$1057 \quad \mathbb{E}[\boldsymbol{\zeta}^t] \leq c'(1 + V(\mathbf{x}))$$

$$1058 \quad \mathbb{E}[\boldsymbol{\eta}^t] \leq c''(1 + V(\mathbf{x})),$$

1059 for all $t = 0, 1, \dots$, where c' and c'' are some positive constants.

1060 C. Mutually dependent observation noise and mutually dependent communication noise.

1061 In this subsection we relax assumptions on observation and communication noises and show that Theorems 4.2 and 4.3 continue to hold. We let Assumptions 1–6 still hold except those which overlap with the following generalizations:

- 1062 • The observation noise \mathbf{n}^t has the joint probability density function p_o such that:

$$1063 \quad \int_{\mathbf{a} \in \mathbb{R}^N} \|\mathbf{a}\| p_o(\mathbf{a}) d\mathbf{a} < \infty, \quad \int_{\mathbf{a} \in \mathbb{R}^N} \mathbf{a} p_o(\mathbf{a}) d\mathbf{a},$$

1064 and $p_o(\mathbf{a}) = p_o(-\mathbf{a})$, for all $\mathbf{a} \in \mathbb{R}^N$.

- 1074 • A (possibly) different nonlinear function $\Psi_{o,i} : \mathbb{R} \rightarrow \mathbb{R}$ is assigned to each
 1075 agent i . Each function $\Psi_{o,i}$ obeys Assumption 3.1.
 1076 • The communication noise ξ_{ij}^t has the joint probability density function $p_{c,ij}$
 1077 such that:

$$1078 \int_{\mathbf{a} \in \mathbb{R}^M} \|\mathbf{a}\| p_{c,ij}(\mathbf{a}) d\mathbf{a} < \infty, \quad \int_{\mathbf{a} \in \mathbb{R}^M} \mathbf{a} p_{c,ij}(\mathbf{a}) d\mathbf{a} = 0,$$

1079 and $p_{c,ij}(\mathbf{a}) = p_{c,ij}(-\mathbf{a})$, for all $\mathbf{a} \in \mathbb{R}^M$.

- 1081 • A different nonlinear function $\Psi_{c,ij,\ell} : \mathbb{R} \rightarrow \mathbb{R}$ is assigned to each arc $(i, j) \in$
 1082 E_d and to each element $\ell = 1, \dots, M$ of the communication noise $[\xi_{ij}^t]_\ell$. Each
 1083 function $\Psi_{c,ij,\ell}$ obeys Assumption 3.1.

1084 This means that observation noises of agents i and j can be mutually dependent.
 1085 Moreover, the communication noises ξ_{ij}^t may have mutually dependent elements $[\xi_{ij}^t]_\ell$,
 1086 for $\ell = 1, \dots, M$. Further, here, for simplicity, we assume that observation and com-
 1087 munication noises are mutually independent.

1088 Let us define functions $\varphi_{o,i} : \mathbb{R} \rightarrow \mathbb{R}$ for $i = 1, 2, \dots, N$ and $\varphi_{ij,\ell} : \mathbb{R} \rightarrow \mathbb{R}$ for $(i, j) \in E$
 1089 and $\ell = 1, 2, \dots, M$ in the same manner as in (4.1), i.e.,

$$1090 (6.9) \quad \varphi_{o,i}(a) = \int \Psi_{o,i}(a+w) p_{o,i}(w) dw,$$

$$1091 (6.10) \quad \varphi_{c,ij,\ell}(a) = \int \Psi_{c,ij,\ell}(a+w) p_{c,ij,\ell}(w) dw.$$

1092 Here, $p_{o,i}$ and $p_{c,ij,\ell}$ are the marginal probability density functions of random variables
 1093 \mathbf{n}_i^t and $[\xi_{ij}^t]_\ell$, respectively. Following same steps as in the proofs of Theorems 4.2
 1094 and 4.3, almost sure convergence and asymptotic normality can be shown. In the
 1095 following, we emphasize only differences. First of all, algorithm (4.4) gets replaced by
 1096

$$1097 \mathbf{x}^{t+1} = \mathbf{x}^t - \alpha_t \left(\frac{b}{a} \hat{\mathbf{L}}_{\varphi_c}(\mathbf{x}^t) - \mathbf{H}^\top \varphi_o(\mathbf{H}((\mathbf{1}_N \otimes \boldsymbol{\theta}^*) - \mathbf{x}^t)) - \mathbf{H}^\top \boldsymbol{\zeta}^t + \frac{b}{a} \boldsymbol{\eta}^t \right).$$

1098 Now, the map $\hat{\mathbf{L}}_{\varphi_c} : \mathbb{R}^{MN} \rightarrow \mathbb{R}^{MN}$ is

$$1100 \hat{\mathbf{L}}_{\varphi_c}(\mathbf{x}) = \begin{bmatrix} \vdots \\ \sum_{j \in \Omega_i} \varphi_{c,ij}(\mathbf{x}_i - \mathbf{x}_j) \\ \vdots \end{bmatrix},$$

1101 for any $\mathbf{x} \in \mathbb{R}^{MN}$, where for all $(i, j) \in E$, function $\varphi_{c,ij} : \mathbb{R}^M \rightarrow \mathbb{R}^M$ is given
 1102 with $\varphi_{c,ij}(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M) = [\varphi_{c,ij,1}(\mathbf{y}_1), \varphi_{c,ij,2}(\mathbf{y}_2), \dots, \varphi_{c,ij,M}(\mathbf{y}_M)]^\top$, for $\mathbf{y} \in \mathbb{R}^M$,
 1103 functions $\varphi_{c,ij,\ell}(a)$ for $(i, j) \in E$ and $\ell = 1, 2, \dots, M$ are given by (6.10). More-
 1104 over, for $\mathbf{y} \in \mathbb{R}^N$, the map $\varphi_o : \mathbb{R}^N \rightarrow \mathbb{R}^N$ is now given with $\varphi_o(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N) =$
 1105 $[\varphi_{o,1}(\mathbf{y}_1), \varphi_{o,2}(\mathbf{y}_2), \dots, \varphi_{o,N}(\mathbf{y}_N)]^\top$. Using the same notation, sequences $\boldsymbol{\zeta}^t \in \mathbb{R}^N$ and
 1106 $\boldsymbol{\eta}^t \in \mathbb{R}^{MN}$ are appropriate versions of the sequences defined in (4.3). If we define
 1107 quantities $\hat{\mathbf{r}}(\mathbf{x})$ and $\hat{\boldsymbol{\gamma}}(t+1, \mathbf{x}, \omega)$ as follows

$$1108 (6.11) \quad \hat{\mathbf{r}}(\mathbf{x}) = -\frac{b}{a} \hat{\mathbf{L}}_{\varphi_c}(\mathbf{x}) - \mathbf{H}^\top \varphi_o(\mathbf{H}(\mathbf{x} - (\mathbf{1}_N \otimes \boldsymbol{\theta}^*))),$$

$$1110 (6.12) \quad \hat{\boldsymbol{\gamma}}(t+1, \mathbf{x}, \omega) = -\frac{b}{a} \boldsymbol{\eta}^t + \mathbf{H}^\top \boldsymbol{\zeta}^t,$$

1111 it is easy to see that all conditions B1–B5 and C1–C5 from Theorem 6.1 still hold
 1112 (see [15]). The only difference occurs in the asymptotic covariance matrix \mathbf{S} , i.e., in
 1113 \mathbf{S}_0 , which is now given by

$$1115 \mathbf{S}_0 = \frac{b^2}{a^2} \mathbf{K}_\eta + \mathbf{H}^\top \mathbf{K}_\zeta \mathbf{H},$$

1116

1117 where $\mathbf{K}_\eta \in \mathbb{R}^{N \times N}$ and $\mathbf{K}_\zeta \in \mathbb{R}^{MN \times MN}$ are the effective covariance matrices of com-
 1118 munication and observation noises after passing through the appropriate nonlinearities
 1119 (analogously defined as cross-covariance matrix $\mathbf{K}_{c,o}$ in Theorem 4.3).

1120 **D. Heavy-tailed noise and identity function.** In this subsection, we show
 1121 that the algorithm (3.1) does not converge in the presence of heavy-tailed observation
 1122 and communication noise if at least one of the nonlinearities Ψ_o and Ψ_c is the identity
 1123 function. This means that in the presence of heavy-tailed observation and communi-
 1124 cation noises, the algorithms from [15, 18] do not converge, in fact, they exhibit an
 1125 infinite variance solution sequence.

1126 **THEOREM 6.3** (Infinite variance). *For the sequence of iterates $\{\mathbf{x}^t\}, t = 1, 2, \dots$,
 1127 generated by (3.1), we have that $\mathbb{E}[\|\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2] = \infty, t = 1, 2, \dots$, if at least one
 1128 of the following statements is true.*

- 1129 1. Function Ψ_o is the identity function, i.e., $\Psi_o(a) = a$ and the observation
 1130 noise has infinite variance, i.e., $\int a^2 d\Phi_o = +\infty$.
- 1131 2. Function Ψ_c is the identity function, i.e., $\Psi_c(a) = a$ and the communication
 1132 noise has infinite variance, i.e., $\int a^2 d\Phi_c = +\infty$.

1133 *Proof.* For simplicity, we assume that if statement 1 holds there is no communi-
 1134 cation noise, i.e. $\boldsymbol{\xi}_{ij} \equiv 0$ for all $(i, j) \in E_d$, and *vice versa*, if statement 2 holds we
 1135 assume that there is no observation noise, i.e., $\mathbf{n} \equiv 0$. If statement 1 holds, in the
 1136 absence of communication noise, the algorithm (3.2) can be written as

$$1137 \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \alpha_t \left(\frac{b}{a} \mathbf{L}_{\Psi_c}(\mathbf{x}) - \mathbf{H}^\top (\mathbf{z}^t - \mathbf{H}\mathbf{x}^t) \right)$$

$$1138 \quad = \mathbf{x}^t - \alpha_t \left(\frac{b}{a} \mathbf{L}_{\Psi_c}(\mathbf{x}) - \mathbf{H}^\top (\mathbf{H}(\mathbf{1} \otimes \boldsymbol{\theta}^*) + \mathbf{n}^t - \mathbf{H}\mathbf{x}^t) \right).$$

1139 If we define $\mathbf{e}^t = \mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*$, $t = 1, 2, \dots$, we have that $\mathbf{e}^{t+1} = \mathbf{F}^t(\mathbf{e}^t) + \alpha_t \mathbf{H}^\top \mathbf{n}^t$,
 1140 where function $\mathbf{F}^t : \mathbb{R}^{MN} \rightarrow \mathbb{R}^{MN}$ is given by $\mathbf{F}^t(\mathbf{y}) = (\mathbf{I} + \alpha_t \mathbf{H}^\top \mathbf{H})\mathbf{y} - \alpha_t \frac{b}{a} \mathbf{L}_{\Psi_c}(\mathbf{y})$,
 1141 for $\mathbf{y} \in \mathbb{R}^{MN}$. Therefore, we have that

$$1142 \quad \|\mathbf{e}^{t+1}\|^2 = \|\mathbf{F}^t(\mathbf{e}^t)\|^2 + 2\alpha_t (\mathbf{F}^t(\mathbf{e}^t))^\top \mathbf{H}^\top \mathbf{n}^t + \alpha_t^2 \|\mathbf{H}^\top \mathbf{n}^t\|^2$$

$$1143 \quad \geq 2\alpha_t (\mathbf{H} \mathbf{F}^t(\mathbf{e}^t))^\top \mathbf{n}^t + \alpha_t^2 \|\mathbf{H}^\top \mathbf{n}^t\|^2,$$

1144 and using the fact that \mathbf{e}^t and \mathbf{n}^t are independent, we have that

$$1145 \quad \mathbb{E}[\|\mathbf{e}^{t+1}\|^2] \geq \alpha_t^2 \mathbb{E}[\|\mathbf{H}^\top \mathbf{n}^t\|^2] = \infty,$$

1146 which completes the proof of statement 1. Proof of statement 2 follows directly from
 1147 Appendix B in [15]. \square

1151 **E. Hypothetical variant of algorithm from [1].** Firstly, we give an overview
 1152 of algorithm that is proposed in [1], for more information see [1]. They considered
 1153 a network of N agents where each agent $i = 1, 2, \dots, N$ at each time $t \geq 0$ collects a
 1154 linear transformation of unknown vector parameter $\mathbf{w}^0 \in \mathbb{R}^M$ corrupted by noise as
 1155 follows

$$1156 \quad d_i(t) = \mathbf{u}_{i,t} \mathbf{w}^0 + v_i(t),$$

1157 where $\mathbf{u}_{i,t} \in \mathbb{R}^M$ is a row regression vector and $v_i(t) \in \mathbb{R}$ is wide-sense stationary zero-
 1158 mean impulsive noise process with variance $\sigma_{v,i}^2$. They introduced an agent-dependent
 1159 and time-varying error nonlinearity, $h_{i,t}(e_i(t))$, into the adaptation step and proposed
 1160 following algorithm

$$1161 \quad \boldsymbol{\psi}_{i,t} = \mathbf{w}_{i,t-1} + \mu_i \mathbf{u}_{i,t}^\top h_{i,t}(e_i(t)),$$

$$1162 \quad (6.13) \quad \mathbf{w}_{i,t} = \sum_{\ell \in \mathcal{N}_i} a_{\ell i} \boldsymbol{\psi}_{\ell,t},$$

1163 where μ_i is a step size parameter, \mathcal{N}_i is the set of agents connected to agent i including
 1164 himself and a_{li} are weighting coefficients. For the error nonlinearity $h_{i,t}(e_i(t))$, they
 1165 set to be a linear combination of $B_i \geq 1$ preselected sign-preserving basis functions,
 1166 i.e., $h_{i,t}(e_i(t)) = \boldsymbol{\alpha}_{i,t}^\top \boldsymbol{\varphi}_{i,t}(e_i(t))$. As it is said in [1], if agent i were to run the sand-
 1167 alone counterpart of the adaptive filter in (6.13), then the optimal nonlinearity that
 1168 minimizes i -th agent MSE is given by $h_{i,t}^{\text{opt}}(x) = -\frac{p'_e(x)}{p_e(x)}$ in terms of the pdf of the
 1169 error signal.

1170 Even though the pdf is not available in practice, for the purpose of comparing algo-
 1171 rithms in the specific numerical example when we know pdf, we introduce hypothetical
 1172 variant of algorithm, by finding optimal $\boldsymbol{\alpha}_{i,t}^{\text{opt}}$, for each agent i at each time t , i.e.,
 1173 $\boldsymbol{\alpha}_{i,t}^{\text{opt}} = \underset{\boldsymbol{\alpha}_{i,t}}{\text{argmin}} \mathbb{E}[h_{i,t}^{\text{opt}}(e_i(t)) - h_{i,t}(e_i(t))]^2$

1174 **F. Derivations and numerical illustrations for Example 1.** Derivation for
 1175 the average per-agent asymptotic variance $\sigma_B^2 = \frac{1}{N} \text{Tr}(\mathbf{S})$ follows

$$1176 \quad \sigma_B^2 = \frac{1}{N} \text{Tr}(a^2 \int_0^{+\infty} e^{\Sigma v} \mathbf{S}_0 e^{\Sigma v} dv) = \frac{1}{N} a^2 \sigma_0^2 h^2 \int_0^{+\infty} \text{Tr}(e^{2\Sigma v} dv)$$

$$1177 \quad = \frac{1}{N} a^2 \sigma_0^2 h^2 \int_0^{+\infty} N e^{(1-2ah^2\varphi'_0(0))v} dv = \frac{a^2 \sigma_0^2 h^2}{2ah^2\varphi'_0(0) - 1}.$$

1178 Integral in the last equality converge for $a > \frac{1}{2h^2\varphi'_0(0)}$.

1180 If $a = a(B) = \frac{1}{2h^2\varphi'_0(0)(B)} + \epsilon$, for some constant $\epsilon > 0$, we have that

$$1181 \quad \sigma_B^2 = \frac{\left(\frac{1}{2h^2\varphi'_0(0)} + \epsilon\right)^2 \sigma_0^2 h^2}{2\left(\frac{1}{2h^2\varphi'_0(0)} + \epsilon\right) h^2 \varphi'_0(0) - 1} = \frac{\left(\frac{1+2h^2\varphi'_0(0)\epsilon}{2h^2\varphi'_0(0)}\right)^2 \sigma_0^2 h^2}{2\left(\frac{1}{2h^2\varphi'_0(0)} + \epsilon\right) h^2 \varphi'_0(0) - 1}$$

$$1182 \quad = \frac{\left(\frac{1+2h^2\varphi'_0(0)\epsilon}{2h^2\varphi'_0(0)}\right)^2 \sigma_0^2 h^2}{1 + 2h^2\varphi'_0(0)\epsilon - 1} = \frac{(1 + 2h^2\varphi'_0(0)\epsilon)^2 \sigma_0^2}{8h^4\varphi'_0(0)^3\epsilon}.$$

1184 Next, we validate that $\lim_{B \rightarrow 0^+} \sigma_B^2 = +\infty$. It is suffice to show that $\lim_{B \rightarrow 0^+} \frac{\sigma_0^2}{\varphi'_0(0)^3} = +\infty$,

1185 since $\sigma_B^2 = \frac{\sigma_0^2}{8h^4\varphi'_0(0)^3\epsilon} + \frac{4h^2\epsilon\sigma_0^2}{8h^4\varphi'_0(0)^2\epsilon} + \frac{4h^4\epsilon^2\sigma_0^2}{8h^4\varphi'_0(0)\epsilon}$.

$$1186 \quad \lim_{B \rightarrow 0^+} \frac{\sigma_0^2}{\varphi'_0(0)^3} = \lim_{B \rightarrow 0^+} \frac{B^2 \int_{-\infty}^{+\infty} \tanh^2\left(\frac{w}{B}\right) f(w) dw}{\left(\int_{-\infty}^{+\infty} \frac{1}{\cosh^2\left(\frac{w}{B}\right)} f(w) dw\right)^3} = \left[\frac{w}{B} = t, dw = dt\right]$$

$$1187 \quad = \lim_{B \rightarrow 0^+} \frac{B^2 \int_{-\infty}^{+\infty} \tanh^2\left(\frac{w}{B}\right) f(w) dw}{B^3 \left(\int_{-\infty}^{+\infty} \frac{1}{\cosh^2(w)} f(Bw) dw\right)^3}$$

$$1188 \quad = \lim_{B \rightarrow 0^+} \frac{\int_{-\infty}^{+\infty} \tanh^2\left(\frac{w}{B}\right) f(w) dw}{B \left(\int_{-\infty}^{+\infty} \frac{1}{\cosh^2(w)} f(Bw) dw\right)^3} = +\infty,$$

1189

1190 since $\lim_{B \rightarrow 0^+} \int_{-\infty}^{+\infty} \tanh^2(\frac{w}{B}) f(w) dw = 1$ and $\lim_{B \rightarrow 0^+} \int_{-\infty}^{+\infty} \frac{1}{\cosh^2(w)} f(Bw) dw < +\infty$.

1191 We now prove that both of the functions σ_o^2 and $\varphi_o'(0)$ are increasing function with
1192 respect to B . Suppose that $B_1 < B_2$, then we have that

$$1193 \quad B_1^2 \tanh^2\left(\frac{w}{B_1}\right) < B_2^2 \tanh^2\left(\frac{w}{B_2}\right),$$

$$1194 \quad \frac{1}{\cosh^2\left(\frac{w}{B_1}\right)} < \frac{1}{\cosh^2\left(\frac{w}{B_2}\right)},$$

1195 for all $w \in \mathbb{R}$. Moreover, since $f(w) \geq 0$ for all $w \in \mathbb{R}$, we have that

$$1197 \quad B_1^2 \tanh^2\left(\frac{w}{B_1}\right) f(w) < B_2^2 \tanh^2\left(\frac{w}{B_2}\right) f(w),$$

$$1198 \quad \frac{1}{\cosh^2\left(\frac{w}{B_1}\right)} f(w) < \frac{1}{\cosh^2\left(\frac{w}{B_2}\right)} f(w),$$

1199 for all $w \in \mathbb{R}$. Therefore, we have that

$$1201 \quad \sigma_o^2(B_1) = \int_{-\infty}^{+\infty} B_1^2 \tanh^2\left(\frac{w}{B_1}\right) f(w) dw < \int_{-\infty}^{+\infty} B_2^2 \tanh^2\left(\frac{w}{B_2}\right) f(w) dw = \sigma_o^2(B_2),$$

$$1202 \quad \varphi_o'(0)(B_1) = \int_{-\infty}^{+\infty} \frac{1}{\cosh^2\left(\frac{w}{B_1}\right)} f(w) dw < \int_{-\infty}^{+\infty} \frac{1}{\cosh^2\left(\frac{w}{B_2}\right)} f(w) dw = \varphi_o'(0)(B_2).$$

1203 We now compare, in the presence of heavy-tailed observation noise with pdf as in (5.3)
1204 for $\beta = 2.05$, the proposed algorithm (5.2) for the optimal choice of B^* with the
1205 method from [1] and its hypothetical variant (see Appendix E). For those methods we
1206 set that $B_i = 2$, $\phi_{i,1}(x) = x$ and $\phi_{i,2}(x) = \tanh(x)$ for all agents. Furthermore, we set
1207 that weighting coefficients are chosen according to $a_{ij} = \frac{\tilde{\mathbf{A}}_{ij}}{\sum_{\ell \in \mathcal{N}_i} \tilde{\mathbf{A}}_{\ell i}}$, where $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$.

1209 Moreover, for the smoothing recursions, zero initial conditions are assumed, ν_i is set
1210 to 0.9 for every agent i and $\epsilon = 10^{-2}$.

1211 Figure 5a shows Monte Carlo estimation of MSE for step size $\alpha_t = \frac{0.5}{t+1}$ and the
1212 Figure 5b shows Monte Carlo estimation of MSE for step size $\alpha_t = \frac{1}{t+1}$. As it can
1213 be seen, the hypothetical variant of the method from [1] outperforms the proposed
1214 one in both of the scenarios. However, that is because with the hypothetical variant
1215 of [1] we optimize the choice of the nonlinearity for each agent at each time, whereas
1216 the proposed algorithm (5.2) is optimized only by average per-agent asymptotic vari-
1217 ance. Moreover, we see that the method from [1] is not as robust as the proposed
1218 algorithm (5.2) with respect to the choice of the step size α_t (constant a).

1219 **G. Proof of the assertion in Remark 4.13.** Here, we modify Theorem 3.1
1220 from [8] and make it applicable to probability density functions that satisfy Assump-
1221 tion 2.2. We will show that

$$1222 \quad (6.14) \quad \sup_{p \in \mathcal{P}_{1+\epsilon}^M} \mathbb{P} \left(|\hat{\theta}_t - \theta^*| > \left(\frac{8^{\frac{1}{\epsilon}} M^{\frac{2}{\epsilon}} \ln 2\delta}{t(\ln 2\delta - 1)} \right)^{\frac{\epsilon}{1+\epsilon}} \right) \geq \delta,$$

1223 for any $\theta^* \in \mathbb{R}$, $\delta \in (0, \frac{1}{2})$, where $\mathcal{P}_{1+\epsilon}^M \subseteq \mathcal{P}$ denotes the subclass of all pdfs from
1224 \mathcal{P} such that $1 + \epsilon$ -central moment equals M for $\epsilon \in (0, 1)$. Therefore, using Markov
1225 inequality, we get

$$1227 \quad \sup_{p \in \mathcal{P}_{1+\epsilon}^M} t \mathbb{E}[|\hat{\theta}_t - \theta^*|^2] \geq c_1 t^{\frac{1-\epsilon}{1+\epsilon}},$$

1228

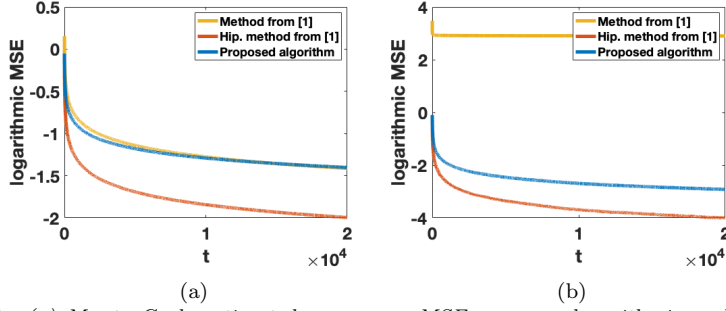


FIG. 5. (a) Monte Carlo-estimated per-sensor MSE error on logarithmic scale for the algorithm (5.2) for optimal B^* and for algorithm and its hypothetical variant from [1] for $a = 0.5$ (b) Monte Carlo-estimated per-sensor MSE error on logarithmic scale for the algorithm (5.2) for optimal B^* and for algorithm and its hypothetical variant from [1] for $a = 1$

1229 for $c_1 = \delta \left(\frac{8^{\frac{1}{\epsilon}} M^{\frac{2}{\epsilon}} \ln 2\delta}{\ln 2\delta - 1} \right)^{\frac{2}{1+\epsilon}}$. Using that $\mathcal{P}_{1+\epsilon}^M \subseteq \mathcal{P}$ and taking the supremum with
 1230 respect to t we get (4.23).

1231 To show that (6.14) holds, we follow the same idea as in [8]. Let us consider the
 1232 class $\mathcal{P}_{+,-} = \{p_+, p_-\}$ of probability density function p_+ and p_- such that p_+ and p_-
 1233 are probability density functions of uniform random variables on $[\frac{p^2-p}{2}, \frac{p^2+p}{2}]$ and on
 1234 $[-\frac{p^2-p}{2}, -\frac{p^2+p}{2}]$, respectively, for $p \in (0, 1)$. It is easy to see that means of probability
 1235 density functions p_+ and p_- are $\theta_+ = \frac{p^2}{2}$ and $\theta_- = -\frac{p^2}{2}$, respectively. Moreover,
 1236 $1 + \epsilon$ -th central moment of both pdfs is equal to

$$1237 \quad (6.15) \quad M = \frac{p^{\epsilon+1}}{2^{\epsilon+1}(\epsilon+2)}.$$

1238 Let $(X_j, Y_j), j = 1, 2, \dots, t$ be i.i.d. pairs random variables such that p_+ is pdf of
 1240 X_1 , and $Y_1 = X_1$ if $X_1 \in I = [\frac{p^2-p}{2}, \frac{p^2+p}{2}]$ and $Y_1 = -X_1$ if $X_1 \notin I$. Notice that
 1241 probability density function of Y_1 is p_- . Since we have that $\mathbb{P}\{X_1 \in I\} = 1 - p$, for
 1242 $X^t = (X_1, X_2, \dots, X_t)$ and $Y^t = (Y_1, Y_2, \dots, Y_t)$, we have that

$$1243 \quad \mathbb{P}\{X^t = Y^t\} = (1-p)^t.$$

1245 Using that $1-p \geq e^{\frac{-p}{1-p}}$, we have that $\mathbb{P}\{X^t = Y^t\} = (1-p)^t \geq 2\delta$, if $p \leq \frac{\ln 2\delta}{\ln 2\delta - t}$.
 1246 Setting that $p := \frac{\ln 2\delta}{t(\ln 2\delta - 1)}$, we have that $p \in (0, 1)$ for all $t = 1, 2, \dots$ and $\delta \in (0, \frac{1}{2})$.

1247 Let $\hat{\theta}_t = \hat{\theta}_t(\cdot)$ be any estimator, then we have that

$$1248 \quad \max \left(\mathbb{P}\left\{ |\hat{\theta}_t(X^t) - \theta_+| > \frac{p^2}{2} \right\}, \mathbb{P}\left\{ |\hat{\theta}_t(Y^t) - \theta_-| > \frac{p^2}{2} \right\} \right)$$

$$1249 \quad \geq \frac{1}{2} \mathbb{P}\left\{ |\hat{\theta}_t(X^t) - \theta_+| > \frac{p^2}{2} \text{ or } |\hat{\theta}_t(Y^t) - \theta_-| > \frac{p^2}{2} \right\}$$

$$1250 \quad \geq \frac{1}{2} \mathbb{P}\{\hat{\theta}_t(X^t) = \hat{\theta}_t(Y^t)\}$$

$$1251 \quad \geq \frac{1}{2} \mathbb{P}\{X^t = Y^t\} \geq \delta.$$

1253 Finally, using (6.15) we get that $\frac{\left(\frac{p^2}{2}\right)^{\frac{\epsilon+1}{2}}}{2\sqrt{2}} \geq \frac{\left(\frac{p^2}{2}\right)^{\frac{\epsilon+1}{2}}}{2^{\frac{\epsilon+1}{2}(\epsilon+1)}} = M \geq Mp^{\frac{\epsilon}{2}}$, which gives us

1254 that $\frac{p^2}{2} \geq \left(8^{\frac{1}{\epsilon}} M^{\frac{2}{\epsilon}} p\right)^{\frac{\epsilon}{\epsilon+1}}$ and therefore we have that

$$1255 \quad \max \left(\mathbb{P} \left\{ |\hat{\theta}_t(X^t) - \theta_+| > \left(\frac{8^{\frac{1}{\epsilon}} M^{\frac{2}{\epsilon}} \ln 2\delta}{t(\ln 2\delta - 1)} \right)^{\frac{\epsilon}{1+\epsilon}} \right\}, \right. \\ 1256 \quad \left. \mathbb{P} \left\{ |\hat{\theta}_t(Y^t) - \theta_-| > \left(\frac{8^{\frac{1}{\epsilon}} M^{\frac{2}{\epsilon}} \ln 2\delta}{t(\ln 2\delta - 1)} \right)^{\frac{\epsilon}{1+\epsilon}} \right\} \right) \geq \delta.$$

1257 Since we have that $\mathcal{P}_{+,-} \subseteq \mathcal{P}_{1+\epsilon}^M$, it follows that (6.14) also holds.

1259 **H. Proof of extensions in Remark 4.14.** For compact notation, we set
1260 that $\bar{\mathbf{H}}$ and $\tilde{\mathbf{H}}^t$ are the $N \times (MN)$ matrices whose i -th row vectors are equal to
1261 $[\mathbf{0}, \dots, \mathbf{0}, (\bar{\mathbf{h}}_i)^\top, \mathbf{0}, \dots, \mathbf{0}]$ and $[\mathbf{0}, \dots, \mathbf{0}, (\tilde{\mathbf{h}}_i^t)^\top, \mathbf{0}, \dots, \mathbf{0}]$, respectively. Hence, for $\mathbf{H}^t =$
1262 $\bar{\mathbf{H}}^t + \tilde{\mathbf{H}}^t$, we have that (4.24) can be written, in compact form, as

$$1263 \quad (6.16) \quad \mathbf{z}^t = \mathbf{H}^t (\mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \mathbf{n}^t = \bar{\mathbf{H}} (\mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \tilde{\mathbf{H}}^t (\mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \mathbf{n}^t.$$

1265 Under this setting, we modify algorithm (3.1) such that, at each time $t = 0, 1, \dots$,
1266 each agent i updates its estimate \mathbf{x}_i^t according to

$$1267 \quad (6.17) \quad \mathbf{x}_i^{t+1} = \mathbf{x}_i^t - \alpha_t \left(\frac{b}{a} \sum_{j \in \Omega_i} \Psi_{\mathbf{c}}(\mathbf{x}_i^t - \mathbf{x}_j^t + \boldsymbol{\xi}_{ij}^t) - \bar{\mathbf{h}}_i \Psi_{\mathbf{o}}(z_i^t - \bar{\mathbf{h}}_i^\top \mathbf{x}_i^t) \right).$$

1268 Assuming that all Assumptions 2.1-3.1 still hold (except those which overlap and
1269 are hence replaced with assumptions in Remark 4.14), we show that the results in
1270 subsections 4.2, 4.3 and 4.4 continue to hold for algorithm (6.17). Following the same
1271 idea as in Section 4, we write algorithm (6.17), in compact form, by:

$$1273 \quad (6.18) \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \alpha_t \left(\frac{b}{a} \mathbf{L}_{\Psi_{\mathbf{c}}}(\mathbf{x}) - \bar{\mathbf{H}}^\top \Psi_{\mathbf{o}}(\mathbf{z}^t - \bar{\mathbf{H}}\mathbf{x}^t) \right).$$

1274 Substituting (6.16) into (6.18), we get that

$$1276 \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \alpha_t \left(\frac{b}{a} \mathbf{L}_{\Psi_{\mathbf{c}}}(\mathbf{x}) - \bar{\mathbf{H}}^\top \Psi_{\mathbf{o}} \left(\bar{\mathbf{H}} (\mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \tilde{\mathbf{H}}^t (\mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \mathbf{n}^t - \bar{\mathbf{H}}\mathbf{x}^t \right) \right) \\ 1277 \quad = \mathbf{x}^t - \alpha_t \left(\frac{b}{a} \mathbf{L}_{\Psi_{\mathbf{c}}}(\mathbf{x}) - \bar{\mathbf{H}}^\top \Psi_{\mathbf{o}} \left(\bar{\mathbf{H}} (\mathbf{1}_N \otimes \boldsymbol{\theta}^* - \mathbf{x}^t) + \tilde{\mathbf{H}}^t (\mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \mathbf{n}^t \right) \right).$$

1278 Recalling $\boldsymbol{\eta}^t \in \mathbb{R}^{MN}$ from (4.3) and defining $\boldsymbol{\zeta}^t \in \mathbb{R}^N$ by

$$1280 \quad \boldsymbol{\zeta}^t = \Psi_{\mathbf{o}} \left(\bar{\mathbf{H}} (\mathbf{1}_N \otimes \boldsymbol{\theta}^* - \mathbf{x}^t) + \tilde{\mathbf{H}}^t (\mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \mathbf{n}^t \right) - \varphi_{\mathbf{o}} \left(\bar{\mathbf{H}} ((\mathbf{1}_N \otimes \boldsymbol{\theta}^*) - \mathbf{x}^t) \right)$$

1281 algorithm (6.18) can be written by

$$1283 \quad (6.19) \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \alpha_t \left(\frac{b}{a} \mathbf{L}_{\varphi_{\mathbf{c}}}(\mathbf{x}^t) - \bar{\mathbf{H}}^\top \varphi_{\mathbf{o}} \left(\bar{\mathbf{H}} ((\mathbf{1}_N \otimes \boldsymbol{\theta}^*) - \mathbf{x}^t) \right) - \bar{\mathbf{H}}^\top \boldsymbol{\zeta}^t + \frac{b}{a} \boldsymbol{\eta}^t \right),$$

1284 Since random variable $\tilde{\mathbf{H}}^t (\mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \mathbf{n}^t$ satisfies Lemma 6.2, the rest of the proofs
1285 are same as in the Section 4.
1286