

1 **DISTRIBUTED RECURSIVE ESTIMATION UNDER HEAVY-TAIL**
2 **COMMUNICATION NOISE**

3 DUSAN JAKOVETIC*, MANOJLO VUKOVIC†, DRAGANA BAJOVIC‡, ANIT KUMAR
4 SAHU§, AND SOUMMYA KAR¶

5 **Abstract.** We consider distributed recursive estimation of an unknown vector parameter
6 $\theta^* \in \mathbb{R}^M$ in the presence of impulsive communication noise. That is, we assume that inter-agent
7 communication is subject to an additive communication noise that may have heavy-tails or is con-
8 taminated with outliers. To combat this effect, within the class of consensus+innovations distributed
9 estimators, we introduce for the first time a nonlinearity in the consensus update. We allow for a
10 general class of nonlinearities that subsumes, e.g., the sign function or component-wise saturation
11 function. For the general nonlinear estimator and a general class of additive communication noises –
12 that may have infinite moments of order higher than one – we establish almost sure (a.s.) convergence
13 to the parameter θ^* . We further prove asymptotic normality and evaluate the corresponding asymp-
14 totic covariance. These results reveal interesting tradeoffs between the negative effect of “loss of
15 information” due to incorporation of the nonlinearity, and the positive effect of communication noise
16 reduction. We also demonstrate and quantify benefits of introducing the nonlinearity in high-noise
17 (low signal-to-noise ratio) and heavy-tail communication noise regimes.

18 **Key words.** Distributed inference; distributed estimation; recursive estimation; heavy-tail
19 noise; consensus+innovations; stochastic approximation.

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

21 **1. Introduction.** We consider distributed inference in networked systems, whe-
22 re each agent in a generic network continuously (over time instances $t = 0, 1, \dots$)
23 makes noisy linear observations of an unknown vector parameter $\theta^* \in \mathbb{R}^M$. Each
24 agent, at each time t , generates a local estimate of θ^* through the so-called consen-
25 sus+innovations strategy, i.e., by 1) weight-averaging its current solution estimate
26 with those of its neighbors, and 2) assimilating its new observation.

27 In this paper, we are interested in consensus+innovations distributed estimation
28 in the presence of an impulsive communication noise, e.g., when the communication
29 noise that corresponds to inter-neighbor communications is heavy-tailed or contami-
30 nated with outliers. It is highly relevant to consider impulsive communication noise
31 in many application scenarios. For example, edge devices in Internet of Things (IoT)
32 systems or sensor networks can be subject to impulsive noise distributions that may
33 not have finite moments of order higher than one, e.g., [8, 32, 13, 37, 12, 9]. In this
34 work, we allow the communication noise to be a zero-mean random variable that
35 may have infinite moments of order α , for any $\alpha > 1$. In particular, communication
36 noise may have an infinite variance. To the best of our knowledge, such scenarios
37 have not been studied in the past work, wherein communication noise in consen-
38 sus+innovations inference is always assumed to have a finite moment of at least second
39 order (finite variance). Actually, as demonstrated ahead in the paper, existing con-

*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 Fundamental Sciences (manojlo.vukovic@uns.ac.rs).

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

§Amazon Alexa AI (anit.sahu@gmail.com).

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

40 sensus+innovation estimators – that are always *linear* in the consensus update part –
 41 can fail to converge under a heavy-tail communication noise. To combat the effect of
 42 the (impulsive or high variance) communication noise, we introduce for the first time
 43 a general nonlinearity in the consensus update. More precisely, we apply a nonlinear
 44 operator (e.g., a sign function, a saturation-like function, or a sigmoid function) on
 45 the difference between an agent’s current iterate and a noisy version of its neighbor’s
 46 iterate, for every agent in the neighborhood set. We establish, under a general setting
 47 for the nonlinearity and the additive communication noise, almost sure (a.s.) conver-
 48 gence of the nonlinear estimator to the true parameter θ^* . We also prove asymptotic
 49 normality and evaluate the corresponding asymptotic covariance in terms of the un-
 50 derlying network topology, observation noise, communication noise, and the employed
 51 nonlinearity. The results reveal interesting interplay among these different problem
 52 dimensions. Most notably, we show that, provided that the nonlinearity has uniformly
 53 bounded outputs, the nonlinear estimator converges a.s. and achieves a finite asymp-
 54 totic covariance, even when the communication noise has no finite moments of order
 55 α for any $\alpha > 1$. We then demonstrate that, in the same regime, the corresponding
 56 linear consensus+innovations estimator has an infinite asymptotic covariance. We
 57 further provide several studies in the finite communication noise variance case that
 58 highlight the regimes where employing the nonlinearity strictly improves performance
 59 of consensus+innovations estimation over linear schemes. Typically, there is a thresh-
 60 old on the communication noise variance above which the nonlinear scheme achieves
 61 a strictly better performance over a linear counterpart.

62 We now review existing literature to help us contrast our contributions with re-
 63 spect to existing work. There has been extensive work on consensus+innovation
 64 distributed estimation, e.g., [17, 15, 16] and related distributed estimation methods,
 65 e.g., [20, 22, 23, 27, 31, 24, 38]. For example, reference [17] derives distributed estima-
 66 tors for both linear and nonlinear observation models, and establishes a.s. convergence
 67 and asymptotic normality of the methods under a general setting for inter-agent com-
 68 munication and observation noises. Specifically, their network model accounts for
 69 random link failures and dithered quantization, which, from the analysis perspective,
 70 effectively translates into an additive communication noise. Reference [15] considers
 71 consensus+innovations distributed estimation in the presence of random link fail-
 72 ures without quantization or additive noise and develops estimators that are asymp-
 73 totically efficient, i.e., that achieve the best achievable asymptotic covariance. The
 74 authors of [16] propose adaptive asymptotically efficient estimators, wherein the inno-
 75 vation gains are adaptively learned during the algorithm progress. There have been
 76 several recent works that consider robust distributed estimation in the presence of
 77 impulsive *observation (sensing)* noise; see [26] for a very recent survey and the refer-
 78 ences therein. To develop robust estimators, various techniques have been utilized,
 79 including, e.g., distributed estimators based on Wilcoxon norm, e.g., [19], Huber loss,
 80 e.g., [21], and mean error minimization, e.g., [36], and novel robust variants of gradi-
 81 ent descent [30]. Reference [1] also considers distributed recursive estimation in the
 82 presence of heavy-tail (impulsive) *sensing (observation) noise* and develops a distrib-
 83 uted estimator that seeks the unknown parameter while at the same time identifying
 84 the optimal error nonlinearity. Reference [6] considers distributed estimation under
 85 measurement attacks. In this setting, the authors develop a consensus+innovations
 86 estimator that employs a saturation nonlinearity in the *innovations update*. Refer-
 87 ences [1, 6] utilize nonlinearities in the *innovations update* to combat the *observation*
 88 *attacks or heavy-tail noise*. This is in contrast with the current paper that employs a
 89 general nonlinearity in the *consensus update* to combat the heavy-tail communication

90 noise. Reference [7] (see also [35]) considers robust distributed estimation methods
 91 based on adaptive subgradient projections. They are also not concerned with com-
 92 bating the effect of heavy-tail inter-agent communication noise. There have also been
 93 several works on consensus+innovations and related distributed detection methods,
 94 e.g., [25, 3, 2, 14]. In particular, reference [14] considers consensus+innovations
 95 distributed detection in the presence of Gaussian additive communication noise. In
 96 summary, with respect to existing work on consensus+innovations distributed infer-
 97 ence, we employ for the first time a general nonlinearity in the consensus update, we
 98 allow for the first time for heavy-tail additive communication noise, and establish for
 99 the considered setting strong convergence guarantees, namely a.s. convergence and
 100 asymptotic normality.

101 The idea of employing a nonlinearity into a “baseline” linear scheme has also been
 102 used in nonlinear versions of the standard average consensus algorithm, e.g., [18, 33, 9].
 103 Average consensus is a distributed algorithm that compute a network-wide average
 104 of scalar values, e.g. [5, 10, 11]. In more detail, the authors of [18] introduce a
 105 trigonometric nonlinearity into a standard linear consensus dynamics and show an
 106 improved dependence of the method on initial conditions. References [33, 9] employ a
 107 general nonlinearity in the linear consensus dynamics and show that it improves the
 108 method’s resilience to additive communication noise. The above works are different
 109 from ours as they focus on the average consensus problem, where the observations are
 110 given to agents beforehand; the corresponding consensus algorithms hence involve only
 111 a consensus step and not an innovation step in the iterative update rule. In contrast,
 112 we consider here the consensus+innovations framework, where new observations are
 113 assimilated at each time instant (algorithm iteration). This technically leads to a
 114 very different analysis with respect to [18, 33, 9], and to qualitatively very different
 115 results. For example, asymptotic performance of the nonlinear consensus+innovations
 116 estimators is determined by an interplay between the effects of network topology,
 117 observation noise and communication noise; observation noise is a model dimension
 118 not present in standard average consensus.

119 There have also been works that employ a specific nonlinearity in the consensus
 120 update within distributed optimization problems. In this context, the authors of [34]
 121 modify the linear consensus update by taking out from the averaging operation the
 122 maximal and minimal estimates among the estimates from all neighbors of an agent.
 123 Reference [4] employs the sign nonlinearity in the consensus update part for distrib-
 124 uted consensus optimization. The works [4, 34] contrast from ours in that they employ
 125 a specific nonlinearity, while we consider a general nonlinearity class. Furthermore,
 126 these works assume deterministic functions in the corresponding distributed consen-
 127 sus optimization problem, that effectively translates into having the observation data
 128 available beforehand. On the other hand, we consider a streaming data scenario that
 129 corresponds to the innovations update part in the algorithm we study.

130 **Paper organization.** Section 2 describes the distributed estimation model that
 131 we consider and presents the nonlinear consensus+innovations estimator that we pro-
 132 pose. Section 3 explains our main results on the almost sure convergence and the
 133 asymptotic normality of the proposed distributed estimator. Section 4 provides sev-
 134 eral analytical and numerical examples that demonstrate benefits of the proposed
 135 nonlinear estimator over the linear counterpart in high and heavy-tail noise regimes.
 136 Finally, Section 5 concludes the paper.

137 **Notation.** We denote by \mathbb{R} the set of real numbers and by \mathbb{R}^m the m -dimensional
 138 Euclidean real coordinate space. We use normal lower-case letters for scalars, lower
 139 case boldface letters for vectors, and upper case boldface letters for matrices. Further,

140 to represent a vector $\mathbf{a} \in \mathbb{R}^m$ through its component, we write $\mathbf{a} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m]^\top$
 141 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
 142 $[\mathbf{A}]_{ij}$, as appropriate, the entry in the i -th row and j -th column of a matrix \mathbf{A} ; \mathbf{A}^\top
 143 the transpose of a matrix \mathbf{A} ; \otimes the Kronecker product of matrices. Further, we use
 144 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 $\mathbf{I}, \mathbf{0}$,
 145 and $\mathbf{1}$ be, respectively, the identity matrix, the zero vector, and the column vector
 146 with unit entries. Further, $\text{Diag}(\mathbf{a})$ is the diagonal matrix whose diagonal entries are
 147 the elements of vector \mathbf{a} ; $\text{Tr}(\mathbf{A})$ the trace of matrix \mathbf{A} ; \mathbf{J} the $N \times N$ matrix $\mathbf{J} :=$
 148 $(1/N)\mathbf{1}\mathbf{1}^\top$. When appropriate, we indicate the matrix or vector dimension through
 149 a subscript. Next, $\mathbf{A} \succ 0$ ($\mathbf{A} \succeq 0$) means that the symmetric matrix A is positive
 150 definite (respectively, positive semi-definite). We further denote by: $\|\cdot\| = \|\cdot\|_2$ the
 151 Euclidean (respectively, spectral) norm of its vector (respectively, matrix) argument;
 152 $\lambda_i(\cdot)$ the i -th smallest eigenvalue; $g'(v)$ the derivative evaluated at v of a function
 153 $g: \mathbb{R} \rightarrow \mathbb{R}$; $\nabla h(\mathbf{w})$ and $\nabla^2 h(\mathbf{w})$ the gradient and Hessian, respectively, evaluated at
 154 w of a function $h: \mathbb{R}^m \rightarrow \mathbb{R}$, $m > 1$; $\mathbb{P}(\mathcal{A})$ and $\mathbb{E}[u]$ the probability of an event \mathcal{A} and
 155 expectation of a random variable u , respectively; and by $\text{sign}(a)$ the sign function,
 156 i.e., $\text{sign}(a) = 1$, for $a > 0$, $\text{sign}(a) = -1$, for $a < 0$, and $\text{sign}(0) = 0$. Finally, for two
 157 positive sequences η_n and χ_n , we have: $\eta_n = O(\chi_n)$ if $\limsup_{n \rightarrow \infty} \frac{\eta_n}{\chi_n} < \infty$.

158 **2. Model and Algorithm.** Subsection 2.1 explains the network and observa-
 159 tion models that we assume. Subsection 2.2 presents the nonlinear consensus+innova-
 160 tions distributed estimator that we propose and states the technical assumptions
 161 needed for subsequent analysis presented in Section 3.

162 **2.1. Problem model.** Consider a network of N agents (sensors). Each agent i
 163 at each time $t = 0, 1, \dots$, collects a linear transformation of the parameter of interest
 164 $\boldsymbol{\theta}^* \in \mathbb{R}^M$, corrupted by noise, as follows:

$$165 \quad (2.1) \quad z_i^t = \mathbf{h}_i^\top \boldsymbol{\theta}^* + n_i^t.$$

166 Here, $z_i^t \in \mathbb{R}$ is the observation, $\mathbf{h}_i \in \mathbb{R}^M$ is the deterministic, non-zero linear trans-
 168 formation vector and $n_i^t \in \mathbb{R}$ is a scalar zero-mean noise. The above update in (2.1)
 169 can be written in a compact form as follows:

$$170 \quad (2.2) \quad \mathbf{z}^t = \mathbf{H}(\mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \mathbf{n}^t.$$

171 Here, $\mathbf{z}^t = [z_1^t, z_2^t, \dots, z_N^t]^\top \in \mathbb{R}^N$ is the observation vector. \mathbf{H} is the $N \times (MN)$
 173 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 i -th
 174 block of size M equals \mathbf{h}_i^\top , and the other M -size blocks are zero vectors; and $\mathbf{n}^t =$
 175 $[n_1^t, n_2^t, \dots, n_N^t]^\top \in \mathbb{R}^N$ is the noise vector at time t .

176 The agents constitute a network $G = (V, E)$, where $V = \{1, \dots, N\}$ is the set of agents,
 177 and E is the set of (undirected) inter-agent communication links (edges) $\{i, j\}$. For
 178 future reference, introduce the $N \times N$ graph Laplacian matrix \mathbf{L} , defined by $\mathbf{L} = \mathbf{D} - \mathbf{A}$,
 179 where \mathbf{D} is the degree matrix and \mathbf{A} is the adjacency matrix. That is, $\mathbf{D} = \text{Diag}(\{d_i\})$,
 180 where d_i is the degree (number of neighbors) of agent i , and \mathbf{A} is a zero-one symmetric
 181 matrix with zero diagonal, such that, for $i \neq j$, $\mathbf{A}_{ij} = 1$ if and only if $\{i, j\} \in E$. Also,
 182 denote by Ω_i the set of neighbors of agent i (excluding i). For an undirected edge
 183 $\{i, j\} \in E$, we denote by (i, j) the arc that points from j to i , and similarly, (j, i) is
 184 the arc that points from i to j . Following this convention, the communication noise
 185 injected when agent j communicates to agent i will be indexed by subscript ij (see
 186 ahead (2.3)).

187 **2.2. Proposed algorithm and technical assumptions.** The agents perform
 188 an iterative consensus+innovations distributed algorithm to collaboratively estimate
 189 the unknown vector parameter $\boldsymbol{\theta}^* \in \mathbb{R}^M$ in the presence of noisy communication links.

190 We assume that communication noise may be heavy-tailed, e.g., [8, 32, 13, 37, 12, 9].
 191 To combat the heavy-tail communication noise, we introduce for the first time a
 192 nonlinear consensus step in consensus+innovations-type methods. More precisely,
 193 the proposed distributed estimator is as follows. At each time $t = 0, 1, \dots$, each agent
 194 i updates its estimate $\mathbf{x}_i^t \in \mathbb{R}^M$ of the parameter $\boldsymbol{\theta}^*$ in the following fashion:

$$195 \quad (2.3) \quad \mathbf{x}_i^{t+1} = \mathbf{x}_i^t - \alpha_t \left(\frac{b}{a} \sum_{j \in \Omega_i} \Psi(\mathbf{x}_i^t - \mathbf{x}_j^t + \boldsymbol{\xi}_{ij}^t) - \mathbf{h}_i(z_i^t - \mathbf{h}_i^\top \mathbf{x}_i^t) \right).$$

196 Here, $\alpha_t = a/(t+1)$ is a step-size, $a, b > 0$ are constants, $\boldsymbol{\xi}_{ij}^t \in \mathbb{R}^M$ is a zero-mean
 197 additive communication noise that models the imperfect communication from agent
 198 j to agent i . Next, $\Psi : \mathbb{R}^M \rightarrow \mathbb{R}^M$ is a non-linear map that operates component-wise
 200 on any vector as follows:

$$200 \quad \Psi(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M) = [\Psi(\mathbf{y}_1), \Psi(\mathbf{y}_2), \dots, \Psi(\mathbf{y}_M)]^\top,$$

203 where, abusing notation, $\Psi : \mathbb{R} \rightarrow \mathbb{R}$ is a component-wise non-linear function. With
 204 algorithm (2.3), upon reception of the noisy version of agent j 's parameter estimate
 205 $\widehat{\mathbf{x}}_{ij}^t = \mathbf{x}_j^t - \boldsymbol{\xi}_{ij}^t$, agent i applies the nonlinearity $\Psi : \mathbb{R}^M \rightarrow \mathbb{R}^M$ on the consensus
 206 contribution $(\mathbf{x}_i^t - \widehat{\mathbf{x}}_{ij}^t)$. Intuitively, the role of Ψ is to combat the communication
 207 noise effect (e.g., truncate large values) while maintaining sufficient useful information
 208 flow. When in algorithm (2.3) we set $\Psi : \mathbb{R}^M \rightarrow \mathbb{R}^M$ to be the identity map, we
 209 recover the \mathcal{LU} (linear estimator) in [17].

210 For future reference, we write algorithm (2.3) in compact form.

211 Let $\mathbf{x}^t = [\mathbf{x}_1^t, \mathbf{x}_2^t, \dots, \mathbf{x}_N^t]^\top \in \mathbb{R}^{MN}$. Furthermore, for $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^\top \in \mathbb{R}^{MN}$ and
 212 $\boldsymbol{\xi} = [\boldsymbol{\xi}_1, \boldsymbol{\xi}_2, \dots, \boldsymbol{\xi}_N]^\top \in \mathbb{R}^{MNN}$, where $\boldsymbol{\xi}_i = [\boldsymbol{\xi}_{i1}, \boldsymbol{\xi}_{i2}, \dots, \boldsymbol{\xi}_{iN}]^\top \in \mathbb{R}^{MN}$ and $\boldsymbol{\xi}_{ij} = 0$ if
 213 $j \notin \Omega_i$, define $\mathbf{L}_\Psi(\mathbf{x}, \boldsymbol{\xi})$ by

$$214 \quad \mathbf{L}_\Psi(\mathbf{x}, \boldsymbol{\xi}) = \begin{bmatrix} \vdots \\ \sum_{j \in \Omega_i} \Psi(\mathbf{x}_i - \mathbf{x}_j + \boldsymbol{\xi}_{ij}) \\ \vdots \end{bmatrix}.$$

215 That is, the map $\mathbf{L}_\Psi(\mathbf{x}, \boldsymbol{\xi}) : \mathbb{R}^{MN} \times \mathbb{R}^{MNN} \rightarrow \mathbb{R}^{MN}$ stacks the N vectors of size M ,
 217 $\sum_{j \in \Omega_i} \Psi(\mathbf{x}_i - \mathbf{x}_j + \boldsymbol{\xi}_{ij})$, $i = 1, 2, \dots, N$, one on top of another. Then, algorithm (2.3)
 218 can be written as:

$$219 \quad (2.4) \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \alpha_t \left(\frac{b}{a} \mathbf{L}_\Psi(\mathbf{x}^t, \boldsymbol{\xi}^t) - \mathbf{H}^\top (\mathbf{z}^t - \mathbf{H}\mathbf{x}^t) \right),$$

220
 221 for $t = 0, 1, \dots$.

222 We make the following assumptions on the underlying network, non-linear map, ob-
 223 servation noise, and communication noise. The assumed nonlinearity class is similar
 224 to that in [29].

225 **Assumption 2.1. Network model:**

226 Graph $G = (V, E)$ is undirected, simple and static.

227 **Assumption 2.2. Nonlinearity Ψ :**

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

- 229 1. Function Ψ is odd, i.e., $\Psi(a) = -\Psi(-a)$, for any $a \in \mathbb{R}$;
- 230 2. $\Psi(a) > 0$, for any $a > 0$;
- 231 3. Function Ψ is a monotonically nondecreasing function;
- 232 4. Ψ is continuous, except possibly on a point set with Lebesgue measure of
 233 zero. Moreover, Ψ is piecewise differentiable.

234 Also, $\Psi : \mathbb{R} \rightarrow \mathbb{R}$ satisfies one of the following two properties:

- 235 5. $|\Psi(a)| \leq c_1(1 + |a|)$, for any $a \in \mathbb{R}$, for some constant $c_1 > 0$;
 236 5'. $|\Psi(a)| \leq c_2$, for some constant $c_2 > 0$.

237 There are many interesting examples of nonlinearities that satisfy Assumption 2.2,
 238 including, e.g., the following:

- 239 • **(NL1)** Sign function: $\Psi(a) = \text{sign}(a)$;
- 240 • **(NL2)** Saturation or clipping function: $\Psi(a) = a$, for $|a| \leq m$; and $\Psi(a) =$
 241 $m \text{sign}(a)$, for $|a| > m$, for some constant $m > 0$;
- 242 • **(NL3)** Relay function with insensitivity zone: $\Psi(a) = 0$, for $|a| \leq r$; and
 243 $\Psi(a) = \text{sign}(a)$, for $|a| > r$, for some constant $r > 0$.

244 **Assumption 2.3. Observation model:**

- 245 1. For each agent $i = 1, \dots, N$, the observation noise sequence $\{n_i^t\}$ in (2.1), is
 246 zero-mean and independent identically distributed (i.i.d.);
- 247 2. Random variables n_i^t and n_j^s are mutually independent whenever the tuple
 248 (i, t) is different from (j, s) ;
- 249 3. Random variable n_i^t has a finite variance equal to σ_{obs}^2 , for any $t = 0, 1, \dots$
 250 and for any $i = 1, \dots, N$;
- 251 4. The matrix $\sum_{i=1}^N \mathbf{h}_i \mathbf{h}_i^\top$ is invertible.

252 The condition 4 in Assumption 2.3 is a standard global observability assumption,
 253 see. e.g. [17]; if it does not hold, then a central estimator that collects all observations
 254 according to (2.1) for each $t = 0, 1, \dots$ and for each $i = 1, \dots, N$, is not able to provide
 255 a consistent sequence of estimates over times $t = 0, 1, \dots$

256 **Assumption 2.4. Communication noise:**

- 257 1. Additive communication noise $\{\xi_{ij}^t\}$, $\xi_{ij}^t \in \mathbb{R}^M$ in (2.3), is i.i.d. in time t ,
 258 independent of the observation noise family $\{n_i^t\}$, $i = 1, \dots, N$, $t = 0, 1, \dots$, and
 259 independent across different arcs (i, j) of graph G ;
- 260 2. Each random variable $[\xi_{ij}^t]_\ell$, for each $t = 0, 1, \dots$, for each arc (i, j) , for each
 261 entry $\ell = 1, \dots, M$, has the same cumulative distribution function Φ ;
- 262 3. The distribution function Φ is symmetric, i.e., for all $a \in \mathbb{R}$ we have that
 263 $\Phi(a) = 1 - \Phi(-a)$, and has strictly positive second moment.

264 We assume that at least one of the conditions 4. or 4'. below holds.

- 265 4. Function Ψ is strictly increasing (from Assumption 2.2) and functions Φ and
 266 Ψ have a common growth point, i.e.,

$$267 \Psi(a_0 + \varepsilon) \geq \Psi(a_0 - \varepsilon),$$

$$268 [\Phi_{ij}]_\ell(a_0 + \varepsilon) \geq [\Phi_{ij}]_\ell(a_0 - \varepsilon),$$

270 for some $a_0 \in \mathbb{R}$ and all $\varepsilon > 0$;

- 271 4'. Distribution Φ has a pdf $p(u)$, $p: \mathbb{R} \rightarrow \mathbb{R}$, that is strictly unimodal, i.e., there
 272 holds $p(0) < +\infty$ and $p(u_1) < p(u_2)$ for $|u_1| > |u_2|$;
- 273 5. There holds that $\int |a| d\Phi(a) < \infty$, and the communication noise is zero-mean,
 274 i.e., $\int a d\Phi(a) = 0$;
- 275 6. If part 5 of Assumption 2.2 holds, then we additionally require that commu-
 276 nication noise has a finite variance, i.e.:

$$277 \int a^2 d\Phi(a) < \infty;$$

- 278 7. Distribution Φ has a well-defined pdf $p: \mathbb{R} \rightarrow \mathbb{R}$ in the vicinity of discontinuity
 279 points of function $\Psi: \mathbb{R} \rightarrow \mathbb{R}$ from Assumption 2.2.

280 For notational simplicity and a clearer presentation, we assume that the com-
 281 munication noise has the same distribution Φ across all arcs (i, j) such that $\{i, j\} \in$
 282 E . We additionally assume that each element of communication noise vector $[\xi_{ij}^t]_\ell$,

283 $\ell = 1, 2, \dots, M$, has the same cumulative distribution function Ψ , and that $[\xi_{ij}^t]_\ell$ and
 284 $[\xi_{ij}^t]_s$ are mutually independent for $\ell \neq s$. Extensions to heterogeneous choices of
 285 nonlinearity Ψ across links and heterogeneous communication noises with mutually
 286 dependent $[\xi_{ij}^t]_\ell$ and $[\xi_{ij}^t]_s$ for $\ell \neq s$, are presented in Remark 1 in Section 3.1 (see also
 287 Supplementary material C). Similarly, we assume that the observation noise has the
 288 same variance across all agents i ; analogous extensions to different agents' observation
 289 noise variances can be performed as well.

290 **3. Main results.** Subsection 3.1 states and proves almost sure convergence of
 291 the proposed nonlinear consensus+innovations distributed estimator in (2.3). Sub-
 292 section 3.2 establishes asymptotic normality of the estimator and evaluates the cor-
 293 responding asymptotic variance.

294 **3.1. Almost sure convergence.** We have the following Theorem.

295 **THEOREM 3.1** (Almost sure convergence). *Let Assumptions 2.1-2.4 hold. Then,*
 296 *for each agent $i = 1, \dots, N$, the sequence of iterates $\{\mathbf{x}_i^t\}$ generated by algorithm (2.3)*
 297 *converges almost surely to the true vector parameter $\boldsymbol{\theta}^*$.*

298 Theorem 3.1 establishes, for a nonlinearity Ψ with bounded outputs (e.g., the
 299 nonlinearities NL1-3 introduced in Section 2), almost sure convergence of the pro-
 300 posed algorithm (2.3) under heavy-tail communication noise that may not have finite
 301 moments of order greater than one. In contrast, it can be shown that the correspond-
 302 ing linear \mathcal{LU} scheme in [17] (obtained by taking Ψ to be the identity function in (2.3))
 303 generates a sequence of iterates with unbounded second moments for all $t = 1, 2, \dots$
 304 (see Supplementary material B). The Theorem also establishes almost sure conver-
 305 gence of (2.3) for nonlinearities with unbounded outputs, more precisely, those that
 306 satisfy part 5 of Assumption 2.2, when the communication noise has finite second
 307 moment. As a special case, by taking Ψ to be the identity map, we recover for the
 308 letter case almost sure convergence of the linear estimator (the \mathcal{LU} algorithm) in [17].

309 **Setting up the proof.** We next outline our strategy for proving Theorem 3.1.
 310 We base our analysis on stochastic approximation arguments. More precisely, we
 311 use Theorem 29 in [17] adapted from [28] (see also Theorem 3 in the supplementary
 312 material) to establish a.s. convergence of \mathbf{x}^t to $\mathbf{1}_N \otimes \boldsymbol{\theta}^*$ by verifying assumptions
 313 B1–B5 of Theorem 29 in [17].

314 The proof strategy is as follows. We first prove a.s. convergence of algorithm (2.3)
 315 for the case without communication noise, i.e., by setting $\xi_{ij}^t \equiv 0$ in (2.3). In this
 316 setting, we first prove the result assuming a continuous function $\Psi : \mathbb{R} \mapsto \mathbb{R}$. Then, we
 317 handle the case with discontinuous Ψ by additionally assuming that we can associate
 318 to $\Psi : \mathbb{R} \mapsto \mathbb{R}$ a “lower bound” surrogate function $\underline{\Psi} : \mathbb{R} \mapsto \mathbb{R}$ that is *continuous*,
 319 satisfies assumption 2.2, and the following holds:

320 (3.1)
$$|\Psi(a)| \geq |\underline{\Psi}(a)|, \text{ for any } a \in \mathbb{R}.$$

321 This enables us to complete the proof for the noiseless case. To transition to the noisy
 322 communications case, a key argument is to consider an auxiliary function $\varphi : \mathbb{R} \rightarrow \mathbb{R}$,
 323 defined by

324 (3.2)
$$\varphi(a) = \int \Psi(a + w) d\Phi(w).$$

325 Intuitively, $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is a convolution-like transformation of nonlinearity $\Psi : \mathbb{R} \rightarrow \mathbb{R}$,
 326 where the convolution is taken with respect to the communication noise cumulative
 327 distribution function Φ .

329 As we will demonstrate ahead, function $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ in the noisy communications case
 330 effectively plays the role that function $\Psi : \mathbb{R} \rightarrow \mathbb{R}$ has in the noiseless case. Moreover,

331 function φ inherits all the key properties of function Ψ . More precisely, we exploit
 332 the following Lemma in [29] (see Lemmas 1-6 in [29]).

333 LEMMA 3.2 ([29]). Consider function φ in (3.2), where function $\Psi : \mathbb{R} \rightarrow \mathbb{R}$,
 334 satisfies Assumption 2.2. Then, the following holds:

- 335 1. φ is odd;
- 336 2. If $|\Psi(\nu)| \leq c_1$, for any $\nu \in \mathbb{R}$, then $|\varphi(a)| \leq c'_2$, for any $a \in \mathbb{R}$, for some
 337 $c'_1 > 0$;
- 338 3. If $|\Psi(\nu)| \leq c_2(1+|\nu|)$, for any $\nu \in \mathbb{R}$, then $|\varphi(a)| \leq c'_2(1+|a|)$, for any $a \in \mathbb{R}$,
 339 for some $c'_2 > 0$;
- 340 4. $\varphi(a)$ is monotonically nondecreasing;
- 341 5. $\varphi(a) > 0$, for any $a > 0$.
- 342 6. φ is continuous at zero;
- 343 7. φ is differentiable at zero, with a strictly positive derivative at zero, equal to:

$$344 \quad (3.3) \quad \varphi'(0) = \sum_{i=1}^s (\Psi(\nu_i + 0) - \Psi(\nu_i - 0)) p(\nu_i) + \sum_{i=0}^s \int_{\nu_i}^{\nu_{i+1}} \Psi'(\nu) p(\nu) d\nu,$$

345 where $\nu_i, i = 1, \dots, s$ are points of discontinuity of Ψ such that $\nu_0 = -\infty$
 346 and $\nu_{s+1} = +\infty$, and we recall that $p(u)$ is the pdf of distribution Φ (see
 347 Assumption 2.2).

348 Lemma 3.2 allows that the treatment of the noisy case becomes completely analogous
 349 to the noiseless case, by replacing function Ψ with φ . Finally, to address the case
 350 when φ may not be continuous over \mathbb{R} , we make use of the following Lemma that is
 351 a trivial corollary of Lemma 3.2.

352 LEMMA 3.3. Consider φ in (3.2). Then, there exists a positive constant ξ such
 353 that $|\varphi(a)| \geq \frac{1}{2}\varphi'(0)|a|$, for $|a| \leq \xi$.

354 Lemma 3.3 allows us to define a continuous function $\underline{\varphi} : \mathbb{R} \mapsto \mathbb{R}$,

$$355 \quad \underline{\varphi}(a) = \begin{cases} \frac{1}{2}\varphi'(0)a & , \quad |a| \leq \xi \\ \xi \text{ sign}(a) & , \quad \text{else} \end{cases},$$

356 that satisfies Assumption 2.2 and obeys the property:

$$358 \quad (3.4) \quad |\varphi(a)| \geq |\underline{\varphi}(a)|, \text{ for any } a \in \mathbb{R}.$$

359 Function $\underline{\varphi}$ will then clearly play the role of function $\underline{\Psi}$ in (3.1) in the noiseless case.

360 We are now ready to prove Theorem 3.1.

361 *Proof.* (Proof of Theorem 3.1)

362 **Step 1: No communication noise.** We start the proof by verifying conditions
 363 B1–B5 of Theorem 29 in [17] for the case without communication noise. We use the
 364 following Lyapunov function $V : \mathbb{R}^{MN} \rightarrow \mathbb{R}$, $V(\mathbf{x}) = \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2$. For this, only
 365 for condition B3, we need to analyze separately the case with continuous Ψ and the
 366 case when Ψ may not be continuous. Also, it can be shown that (2.3) can be put in
 367 the form required by Theorem 29 in [17] (see also (36) in the supplementary material)
 368 by letting

$$369 \quad (3.5) \quad \mathbf{r}(\mathbf{x}) = -\mathbf{H}^\top \mathbf{H}(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) - \frac{b}{a} \mathbf{L}_\Psi(\mathbf{x}, \mathbf{0}),$$

$$370 \quad (3.6) \quad \gamma(t+1, \mathbf{x}, \omega) = \mathbf{H}^\top \mathbf{n}^t,$$

371 where ω denotes an element of the underlying probability space.

372 Consider the filtration $\mathcal{F}_t, t = 1, 2, \dots$, where \mathcal{F}_t is the σ - algebra generated by $\{\mathbf{n}^s\}_{s=0}^{t-1}$.
 373 Denote by $(\Omega, \mathcal{F}, \mathbb{P})$ the underlying probability space that generates random vectors
 374 $\mathbf{n}^t, t = 0, 1, 2, \dots$, and by $\omega \in \Omega$ its arbitrary element. Clearly, for each t , function
 375 $\gamma(t+1, \cdot, \cdot)$ is $\mathcal{B}^{MN} \otimes \mathcal{F}$ measurable, where \mathcal{B}^{MN} is the Borel sigma algebra on \mathbb{R}^{MN} .
 376

377 Also, $\mathbf{r}(\cdot)$ is \mathcal{B}^{MN} measurable. Hence, condition B1 holds. Further, the family of
 378 random vectors $\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)$ is \mathcal{F}_t measurable, zero-mean and independent of \mathcal{F}_{t-1} .
 379 Thus, condition B2 holds.

380 We now inspect condition B3. Assume first that function $\Psi : \mathbb{R} \mapsto \mathbb{R}$ is contin-
 381 uous. The gradient of V equals $\nabla V(\mathbf{x}) = 2(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)$. Clearly, function $V(\cdot)$
 382 is twice continuously differentiable and has uniformly bounded second order partial
 383 derivatives. We consider

$$384 \quad (3.7) \quad S = \sup_{\|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\| \in (\epsilon, 1/\epsilon)} \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle, .$$

385 We will show that $S < 0$, thus verifying condition B3. We have, for any $\mathbf{x} \in \mathbb{R}^{MN}$:

$$387 \quad \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle = -2(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \left(\mathbf{H}^\top \mathbf{H} (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) + \frac{b}{a} \mathbf{L}_\Psi(\mathbf{x}) \right)$$

$$388 \quad (3.8) \quad = -2 \underbrace{\left((\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{H}^\top \mathbf{H} (\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) \right)}_{T_1(\mathbf{x})} - 2 \frac{b}{a} \underbrace{(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{L}_\Psi(\mathbf{x})}_{T_2(\mathbf{x})}.$$

389 Clearly $T_1 = T_1(\mathbf{x}) \geq 0$. We will also show that $T_2 = T_2(\mathbf{x}) \geq 0$. Utilizing the fact
 390 that, $\Psi(\cdot)$ is an odd function, we have that,

$$392 \quad (3.9) \quad T_2 = \sum_{\{i,j\} \in E, i < j} (\mathbf{x}_i - \mathbf{x}_j)^\top \Psi(\mathbf{x}_i - \mathbf{x}_j) \geq 0,$$

393 as for $\mathbf{g} = (\mathbf{x}_i - \mathbf{x}_j)$, we have that,

$$395 \quad (3.10) \quad (\mathbf{x}_i - \mathbf{x}_j)^\top \Psi(\mathbf{x}_i - \mathbf{x}_j) = \sum_{\ell=1}^M \mathbf{g}_\ell^\top \Psi(\mathbf{g}_\ell) \geq 0,$$

396 because \mathbf{g}_ℓ and $\Psi(\mathbf{g}_\ell)$ have the same sign, by Assumption 2.2. Therefore,

$$398 \quad \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle = -2T_1 - 2 \frac{b}{a} T_2 \leq 0,$$

399 for any $\mathbf{x} \in \mathbb{R}^{MN}$.

401 We will further show that S in (3.7) is strictly less than 0. First, consider the set
 402 $\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^{MN} : \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\| \in [\epsilon, 1/\epsilon]\}$. Note that set \mathcal{C} is nonempty and compact.
 403 Clearly, we have that:

$$404 \quad (3.11) \quad S \leq S_{\mathcal{C}} := \sup_{\mathbf{x} \in \mathcal{C}} \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle,$$

405 It is thus sufficient to show that $S_{\mathcal{C}} < 0$. Suppose the contrary is true, i.e., suppose
 406 that $S_{\mathcal{C}} = 0$. As set \mathcal{C} is compact and function $\mathbf{x} \mapsto \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle$ is continuous, by the
 407 Weierstrass theorem, we have that $S_{\mathcal{C}} = 0$ is equivalent to having $\langle \mathbf{r}(\mathbf{x}^\bullet), \nabla V(\mathbf{x}^\bullet) \rangle = 0$,
 408 for some point $\mathbf{x}^\bullet \in \mathcal{C}$. In this case, \mathbf{x}^\bullet has to be of the form, $\mathbf{x}^\bullet = \mathbf{1}_N \otimes \mathbf{m}$, where $\mathbf{m} \in$
 409 \mathbb{R}^M . As otherwise, we would have that, T_2 is strictly positive. But then, we have, $T_1 =$
 410 $\left((\mathbf{x}^\bullet - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)^\top \mathbf{H}^\top \mathbf{H} (\mathbf{x}^\bullet - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) \right) = (\mathbf{m} - \boldsymbol{\theta}^*)^\top \left(\sum_{i=1}^N \mathbf{h}_i \mathbf{h}_i^\top \right) (\mathbf{m} - \boldsymbol{\theta}^*) > 0$,
 411 which is a contradiction in view of (3.7). Hence, we conclude that, for a continuous
 412 function Ψ , it holds that $S < 0$, and that condition B3 holds, i.e.,

$$413 \quad \sup_{\|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\| \in (\epsilon, 1/\epsilon)} \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle < 0.$$

414 Now, we verify condition B3 for function Ψ that is not continuous but to which we can
 415 associate function $\underline{\Psi}$ that obeys Assumption 2.2 and for which condition (3.1) holds.
 416 Then, the verification of condition B3 follows analogously to the case with continuous
 417 Ψ by replacing T_2 in (3.9) with the following lower bound of T_2

$$419 \quad (3.12) \quad \underline{T}_2 = \sum_{\{i,j\} \in E, i < j} (\mathbf{x}_i - \mathbf{x}_j)^\top \underline{\Psi}(\mathbf{x}_i - \mathbf{x}_j),$$

420 where $\underline{\Psi}(\mathbf{a}) = [\underline{\Psi}(\mathbf{a}_1), \underline{\Psi}(\mathbf{a}_2), \dots, \underline{\Psi}(\mathbf{a}_M)]^\top$. Hence, condition B3 is verified.

We next verify condition B4. Recalling the definition of $\mathbf{r}(\mathbf{x})$ in (3.5), we have,

$$(3.13) \quad \|\mathbf{r}(\mathbf{x})\|^2 \leq c_3 V(\mathbf{x}) + c_4 \|\Psi(\mathbf{x})\|^2,$$

where $c_3 = 2a^2 \|\mathbf{H}^\top \mathbf{H}\|^2$ and $c_4 = 2b^2 \|\mathbf{L}\|^2$.

We also have that,

$$\|\Psi(\mathbf{x})\| \leq c_5 \sum_{\{i,j\} \in E} (|\mathbf{x}_i - \boldsymbol{\theta}^*| + |\mathbf{x}_j - \boldsymbol{\theta}^*|) + c_6 \leq c_7 \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\| + c_6,$$

for some positive constants c_5, c_6, c_7 . Therefore, we have

$$(3.14) \quad \|\Psi(\mathbf{x})\|^2 \leq 2c_7 V(\mathbf{x}) + 2c_8^2,$$

for some positive constant c_8 .

Thus, we have that,

$$\|\mathbf{r}(\mathbf{x})\|^2 \leq c_9 V(\mathbf{x}) + c_{10},$$

for some positive constants c_9, c_{10} . Recall $\gamma(t+1, \mathbf{x}, \omega)$ in (3.6). Using the boundedness of the second moment of the observation noise, we finally have that,

$$\|\mathbf{r}(\mathbf{x})\|^2 + \mathbb{E} \left[\|\gamma(t+1, \mathbf{x}^t, \omega)\|^2 \right] \leq c_{11} (V(\mathbf{x}) + 1),$$

for some positive constant c_{11} . Hence, condition B4 is satisfied. Finally, condition B5 clearly holds. Therefore, we conclude that $\mathbf{x}^t \rightarrow \mathbf{1}_N \otimes \boldsymbol{\theta}^*$, almost surely.

Step 2: The case with communication noise. We proceed by considering algorithm (2.3) under communication noise.

We clarify the steps needed to transition from the noiseless to the noisy case. If we write

$$\Psi(\mathbf{x}_i^t - \mathbf{x}_j^t + \boldsymbol{\xi}_{ij}^t) = \boldsymbol{\varphi}(\mathbf{x}_i^t - \mathbf{x}_j^t) + \boldsymbol{\eta}_{ij}^t,$$

where $\boldsymbol{\eta}_{ij}^t = [\Psi(\mathbf{x}_i^t - \mathbf{x}_j^t + \boldsymbol{\xi}_{ij}^t) - \boldsymbol{\varphi}(\mathbf{x}_i^t - \mathbf{x}_j^t)]$ and $\boldsymbol{\varphi} : \mathbb{R}^M \rightarrow \mathbb{R}^M$ is component-wise map defined as $\boldsymbol{\varphi}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M) = [\varphi(\mathbf{x}_1), \varphi(\mathbf{x}_2), \dots, \varphi(\mathbf{x}_M)]^\top$. We will see that quantity $\boldsymbol{\eta}_{ij}^t$ is a key ingredient of $\gamma(t+1, \mathbf{x}, \omega)$ in Theorem 29 in [17] (see also Theorem 3 in the supplementary material).

The algorithm (2.3) can be written in compact form:

$$(3.15) \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \alpha_t \left(\frac{b}{a} \mathbf{L}_\varphi(\mathbf{x}^t) - \mathbf{H}^\top (\mathbf{z}^t - \mathbf{H}\mathbf{x}^t) + \frac{b}{a} \boldsymbol{\eta}^t \right).$$

Here,

$$(3.16) \quad \mathbf{L}_\varphi(\mathbf{x}^t) = \begin{bmatrix} \vdots \\ \sum_{j \in \Omega_i} \boldsymbol{\varphi}(\mathbf{x}_i^t - \mathbf{x}_j^t) \\ \vdots \end{bmatrix} \in \mathbb{R}^{MN}, \quad \boldsymbol{\eta}^t = \begin{bmatrix} \vdots \\ \sum_{j \in \Omega_i} \boldsymbol{\eta}_{ij}^t \\ \vdots \end{bmatrix} \in \mathbb{R}^{MN},$$

where the $M \times 1$ blocks $\sum_{j \in \Omega_i} \boldsymbol{\varphi}(\mathbf{x}_i^t - \mathbf{x}_j^t)$ and $\sum_{j \in \Omega_i} \boldsymbol{\eta}_{ij}^t$ are stacked one on top of another

for $j = 1, \dots, N$.

The differences of (3.15) with respect to the case without additive communication noise are that \mathbf{L}_φ replaces \mathbf{L}_Ψ and the term $\frac{b}{a} \alpha_t \boldsymbol{\eta}^t$ is added.

We define the Lyapunov function $V : \mathbb{R}^{MN} \rightarrow \mathbb{R}$, and quantities $\mathbf{r}_\varphi(x)$ and $\gamma_\varphi(t, \mathbf{x}, \omega)$ as follows:

$$(3.17) \quad V(\mathbf{x}) = \|\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2$$

$$(3.18) \quad \mathbf{r}_\varphi(\mathbf{x}) = -\mathbf{H}^\top \mathbf{H}(\mathbf{x} - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) - \frac{b}{a} \mathbf{L}_\varphi(\mathbf{x}),$$

$$(3.19) \quad \gamma_\varphi(t+1, \mathbf{x}, \omega) = \mathbf{H}^\top \mathbf{n}^t - \frac{b}{a} \boldsymbol{\eta}^t,$$

Now, make the following identification with respect to the transition from the noiseless to the noisy case. Quantity $\mathbf{H}^\top \mathbf{n}^t$ in the noiseless case is replaced with quantity

469 $\mathbf{H}^\top \mathbf{n}^t - \frac{b}{a} \boldsymbol{\eta}^t$ in the noisy case. The map $\mathbf{L}_\Psi(\cdot, \mathbf{0}) : \mathbb{R}^{MN} \rightarrow \mathbb{R}^{MN}$ in (2.4) is replaced
 470 with the map $\mathbf{L}_\varphi : \mathbb{R}^{MN} \rightarrow \mathbb{R}^{MN}$ given in (3.16).

471 The proof proceeds analogously by again verifying Assumptions B1–B5. We only
 472 clarify the differences in verifying these conditions with respect to the noiseless case.
 473 The filtration \mathcal{F}_t is replaced with the filtration \mathcal{G}_t , $t = 1, 2, \dots$, which is generated
 474 not only by $\{\mathbf{n}^s\}_{s=0}^{t-1}$ but also by $\{\boldsymbol{\xi}_{ij}^s\}_{s=0}^{t-1}$ for $(i, j) \in E$. Clearly, for each t , function
 475 $\gamma_\varphi(t+1; \cdot; \cdot)$ is $\mathcal{B}^{MN} \otimes \mathcal{F}$ measurable. Also, $\mathbf{r}_\varphi(\cdot)$ is \mathcal{B}^{MN} measurable. Hence, condition
 476 B1 holds. Further, the family of random vectors $\gamma_\varphi(t+1, \mathbf{x}, \omega)$ is \mathcal{F}_t measurable,
 477 zero-mean and independent of \mathcal{F}_{t-1} . Thus, condition B2 holds. As function φ is
 478 odd, non-decreasing, strictly positive for its positive arguments, and has a positive
 479 derivative at zero by Lemma 3.2, condition B3 is derived analogously to the noiseless
 480 case. Conditions B4 and B5 hold analogously to the noiseless case. Thus, the result
 481 is verified. \square

482 **Remark 1:** Theorem 3.1 continues to hold under the following generalizations:

- 483 • A different nonlinear function $\Psi_{ij,\ell} : \mathbb{R} \rightarrow \mathbb{R}$ is assigned to each arc (i, j) and
 484 to each element $\ell = 1, \dots, M$ of the communication noise $[\boldsymbol{\xi}_{ij}^t]_\ell$. Each function
 485 $\Psi_{ij,\ell}$ obeys Assumption 2.2.
- 486 • The observation noise $\sigma_{obs,i}^2$ is different for each agent $i = 1, 2, \dots, N$.
- 487 • The communication noise $\boldsymbol{\xi}_{ij}^t$ has the joint cumulative distribution function
 488 Φ_{ij} such that:

$$489 \int_{\mathbf{a} \in \mathbb{R}^M} \|\mathbf{a}\| d\Phi_{ij}(\mathbf{a}) < \infty, \quad \int_{\mathbf{a} \in \mathbb{R}^M} \mathbf{a} d\Phi_{ij}(\mathbf{a}) = 0,$$

490 and $\Phi_{ij}(\mathbf{a}) = 1 - \Phi_{ij}(-\mathbf{a})$, for all $\mathbf{a} \in \mathbb{R}^M$.

492 All the remaining assumptions in 2.1-2.4 continue to hold.

493 Note that the above means that the communication noise $\boldsymbol{\xi}_{ij}^t$ may have mutually
 494 dependent elements $[\boldsymbol{\xi}_{ij}^t]_\ell$, for $\ell = 1, \dots, M$.

495 For the above generalization, it can be shown that Theorem 3.1 continues to hold (see
 496 Supplementary material C).

497 **3.2. Asymptotic normality.** We now present our results on asymptotic nor-
 498 mality of estimator (2.3).

499 **THEOREM 3.4 (Asymptotic normality).** *Let Assumptions 2.1 – 2.4 hold. Con-*
 500 *sider algorithm (2.3) with step-size $\alpha_t = a/(t+1)$, $t = 0, 1, \dots$, $a > 0$. Then, the*
 501 *normalized sequence of iterates $\{\sqrt{t+1}(\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*)\}$ converges in distribution to a*
 502 *zero-mean multivariate normal random vector, i.e., the following holds:*

$$503 \sqrt{t+1}(\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*) \Rightarrow \mathcal{N}(\mathbf{0}, \mathbf{S}),$$

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

$$506 (3.20) \quad \mathbf{S} = a^2 \int_0^\infty e^{\boldsymbol{\Sigma}v} \mathbf{S}_0 e^{\boldsymbol{\Sigma}^\top v} dv.$$

507 Here, $\mathbf{S}_0 = \sigma_{obs}^2 \mathbf{H}^\top \mathbf{H} + \frac{b^2}{a^2} \sigma^2 \text{Diag}(\{d_i \mathbf{I}_M\})$, where we recall that d_i is the degree of
 508 agent i ; $\sigma^2 = \int |\Psi(w)|^2 d\Phi(w)$ is the effective communication noise variance after
 509 passing through the nonlinearity Ψ ; we recall the observation matrix \mathbf{H} in (2.2); the
 510 observation noise variance σ_{obs}^2 in (2.1); function ϕ in (3.2); and $\boldsymbol{\Sigma} = \frac{1}{2} \mathbf{I} - a(\mathbf{H}^\top \mathbf{H} +$
 511 $\frac{b}{a} \varphi'(0)(\mathbf{L} \otimes \mathbf{I}_M))$, where a is taken large enough such that matrix $\boldsymbol{\Sigma}$ is stable (i.e., real
 512 parts of $\boldsymbol{\Sigma}$'s eigenvalues are negative).

513 Theorem 3.4 shows that, for the communication noise with finite variance and un-

514 bounded nonlinearities that satisfy part 5 of Assumption 2.2, the variance with the
 515 proposed nonlinear estimator (2.3) decays (in the weak convergence sense) at (the best
 516 achievable) rate $O(1/t)$. In particular, by taking Ψ to be the identity function, we
 517 recover the asymptotic normality result in [17] of the corresponding linear estimator
 518 (the \mathcal{LU} scheme in [17]). Note also that the asymptotic variance expression in (3.20)
 519 for the identity function $\Psi(a) = a$ coincides with that in [17] for the \mathcal{LU} scheme.

520 Theorem 3.4 further demonstrates that, even under a heavy-tailed communication
 521 noise (with unbounded variance) with a bounded nonlinearity (e.g., nonlinearities
 522 NL1-3 in Section 2), the variance with algorithm (2.3) still decays at rate $O(1/t)$. In
 523 contrast, the corresponding linear scheme (obtained by taking Ψ in (2.3) to be the
 524 identity function) generates a sequence with unbounded variances for each $t = 1, 2, \dots$
 525 More precisely, we then have that $E[\|\mathbf{x}^t - \mathbf{1}_N \otimes \boldsymbol{\theta}^*\|^2] = \infty$, for any $t = 1, 2, \dots$ (see
 526 Supplementary material C).

527 Theorem 3.4 explicitly quantifies asymptotic variance of (2.3). This also allows, in
 528 the finite communication noise regime, to compare the nonlinear versus the linear
 529 scheme (when both schemes achieve a finite asymptotic variance). See Subsection 4.1
 530 for details.

531 Theorem 3.4 also reveals an interesting tradeoff when including the nonlinearity Ψ
 532 into the consensus update. On the one hand, nonlinearity makes a beneficial effect in
 533 that the communication noise plays the role only through the effective variance $\sigma^2 =$
 534 $\int |\Psi(w)|^2 d\Phi(w)$. In contrast, with the linear scheme, σ^2 is replaced with $\int w^2 d\Phi(w)$
 535 that is infinite under a heavy tail setting. On the other hand, the nonlinearity Ψ
 536 makes a negative effect in that it “reduces quality” of matrix $\boldsymbol{\Sigma}$ through the quantity
 537 $\varphi'(0)$ that is typically less than one with a nonlinear scheme and equal to one with the
 538 linear scheme. Clearly, the tradeoff goes in favor of the nonlinear scheme in the heavy
 539 tail setting (finite variance with the nonlinear estimator versus infinite variance with
 540 the linear estimator). In the finite communication noise variance setting, the nonlinear
 541 scheme typically improves performance under a sufficiently low communication signal
 542 to noise ratio (SNR); see also Subsection 4.1. We are now ready to prove Theorem
 543 3.4.

544 *Proof.* (Proof of Theorem 3.4) We establish asymptotic normality by verifying
 545 assumptions C1-C5 of Theorem 29 in [17] (see also Theorem 3 in the supplementary
 546 material). Firstly, we show that condition C1 hold. Since function φ is differentiable
 547 at zero, we have that

$$548 \quad (3.21) \quad \varphi(a) = \varphi(0) + \varphi'(0)a + \Delta(a) = \varphi'(0)a + \Delta(a),$$

550 where for the function $\Delta : \mathbb{R} \rightarrow \mathbb{R}$, we have that $\lim_{a \rightarrow 0} \frac{\Delta(a)}{a} = 0$. Hence, the function
 551 $\mathbf{r}_\varphi(\mathbf{x})$ admits representation as in Theorem 29 of [17] (see also (37) of Theorem 3 in
 552 the supplementary material), with matrix

$$553 \quad \mathbf{B} = -\mathbf{H}^T \mathbf{H} - \frac{b}{a} \varphi'(0) [\mathbf{L} \otimes \mathbf{I}_M],$$

554 and function $\boldsymbol{\delta} : \mathbb{R}^{MN} \rightarrow \mathbb{R}^{MN}$, given with $\boldsymbol{\delta}(\mathbf{x}) = -\frac{b}{a} \mathbf{L}_\Delta(\mathbf{x})$. Here, function $\mathbf{L}_\Delta(\mathbf{x}) :$
 555 $\mathbb{R}^{MN} \rightarrow \mathbb{R}^{MN}$ is defined by

$$557 \quad \mathbf{L}_\Delta(\mathbf{x}) = \begin{bmatrix} \vdots \\ \sum_{j \in \Omega_i} \Delta(\mathbf{x}_i - \mathbf{x}_j) \\ \vdots \end{bmatrix},$$

558 where function $\Delta : \mathbb{R}^M \rightarrow \mathbb{R}^M$ is defined by (3.21), $\Delta(\mathbf{y}_1, \mathbf{y}_1, \dots, \mathbf{y}_M) = [\Delta(\mathbf{y}_1), \Delta(\mathbf{y}_2), \dots, \Delta(\mathbf{y}_M)]^\top$,
 559 $\mathbf{y} \in \mathbb{R}^M$. ■

561 Condition C2 trivially holds, if we use that $\alpha_t = \frac{a}{t+1}$. Furthermore, $\Sigma = a\mathbf{B} + \frac{1}{2}\mathbf{I}$ is
 562 stable if a is large enough, because matrix $-\mathbf{B}$ is positive definite (see [17]). Thus,
 563 condition C3 also holds.

564 For $\mathbf{A}(t, \mathbf{x}) = \mathbb{E}[\gamma_\varphi(t+1, \mathbf{x}, \omega)\gamma_\varphi^\top(t+1, \mathbf{x}, \omega)]$, using the Lebesgue's dominated
 565 convergence theorem, it can be shown that

$$566 \lim_{t \rightarrow \infty, \mathbf{x} \rightarrow \theta^*} \mathbf{A}(t, \mathbf{x}) = \sigma_{\text{obs}}^2 \mathbf{H}^\top \mathbf{H} + \sigma^2 \text{Diag}(\{d_i \mathbf{I}_M\}).$$

567 Therefore, condition C4 also holds. It remains to verify condition C5. Recall quantity
 568 $\gamma_\varphi(t+1, \mathbf{x}, \omega)$ in (3.19). Note that this condition is equivalent to saying that the family
 569 of random variables $\{\|\gamma_\varphi(t+1, \mathbf{x}, \omega)\|^2\}_{t=0,1,\dots, \|\mathbf{x}-\theta^*\|<\epsilon}$ is uniformly integrable. If
 570 the condition 5 in Assumption 2.2 holds (the case with finite communication noise
 571 variance and the nonlinearity with unbounded outputs), then:

$$572 (3.22) \quad \|\gamma_\varphi(t+1, x, \omega)\|^2 \leq c_{12} + c_{13}\|\mathbf{n}^t\|^2 + c_{14}\|\boldsymbol{\eta}^t\|^2,$$

573 for some positive constants c_{12}, c_{13}, c_{14} .

574 Consider the family $\{\tilde{\mathbf{g}}(t+1, \mathbf{x}, \omega)\}_{t=0,1,\dots, \|\mathbf{x}-\theta^*\|<\epsilon}$, with

$$575 (3.23) \quad \tilde{\mathbf{g}}(t+1, \mathbf{x}, \omega) = c_{12} + c_{13}\|\mathbf{n}^t\|^2 + c_{14}\|\boldsymbol{\eta}^t\|^2.$$

576 Clearly, $\tilde{\mathbf{g}}(t+1, x, \omega)$ is integrable, for any $t = 0, 1, \dots$, for any $\epsilon > 0$, due to
 577 the finite second moment of sensing and observation noises. The family $\{\tilde{\mathbf{g}}(t+1, x, \omega)\}_{t=0,1,\dots, \|\mathbf{x}-\theta^*\|<\epsilon}$
 578 is i.i.d. and hence it is uniformly integrable. The family
 579 $\{\|\gamma_\varphi(t+1, x, \omega)\|^2\}_{t=0,1,\dots, \|\mathbf{x}-\theta^*\|<\epsilon}$ is dominated by $\{\tilde{\mathbf{g}}(t+1, x, \omega)\}_{t=0,1,\dots, \|\mathbf{x}-\theta^*\|<\epsilon}$
 580 that is uniformly integrable, and hence $\{\|\gamma_\varphi(t+1, x, \omega)\|^2\}_{t=0,1,\dots, \|\mathbf{x}-\theta^*\|<\epsilon}$ is also uni-
 581 formly integrable. An analogous argument can be applied if condition 5' in Assump-
 582 tion 2.2 holds (bounded nonlinearity, communication noise with infinite variance).
 583 Hence, condition C5 holds; thus, the result. \square

584 **4. Analytical and numerical examples.** Subsection 4.1 provides analytical
 585 examples, and Subsection 4.2 provides simulation examples, that illustrate the main
 586 results presented in Section 3.

587 **4.1. Analytical examples.** We provide several analytical examples that illus-
 588 trate Theorem 3.4. The examples demonstrate that, in the considered setting, the
 589 proposed nonlinear method in (2.3) achieves a lower asymptotic variance than the
 590 corresponding linear scheme, for a low SNR regime, i.e., for the case when the com-
 591 munication noise variance is above a threshold. We also consider optimization of the
 592 nonlinearity Ψ for a given nonlinearity class; more precisely, for the given analytical
 593 example, we consider optimization of parameter B for the NL2 nonlinearity class in
 594 Section 2.

595 **Example 1:** We follow a setup similar to [17], but we consider the nonlinear con-
 596 sensus+innovations scheme in (2.3), with the non-linear operator $\Psi : \mathbb{R} \rightarrow \mathbb{R}$ of the
 597 following form (the NL2 nonlinearity):

$$598 (4.1) \quad \Psi(w) = \begin{cases} w & , \quad |w| \leq B \\ +B & , \quad w > B \\ -B & , \quad w < B \end{cases},$$

599 for some parameter $B > 0$. Notice that letting $B \rightarrow \infty$ in (4.1) leads to the linear
 600 consensus+innovations \mathcal{LU} scheme in [17].

601 Each agent i observes a scalar parameter $\theta^* \in \mathbb{R}$ according to:

$$602 z_i(t) = h\theta^* + n_i^t,$$

603 where $h \neq 0$ and n_i^t is i.i.d. in time and across sensors with variance σ_{obs}^2 and zero
 604 mean. Communication noise is i.i.d. across arcs and in time and is independent of
 605 $\{n_i^t\}$, for all $i = 1, 2, \dots, N$. Assume that the communication noise has a probability

609 distribution function $f(w)$ that is strictly positive in the vicinity of zero. Denote the
 610 eigenvalues of \mathbf{L} by $0 = \lambda_1 < \lambda_2 \leq \dots \leq \lambda_N$. Let the graph be regular, for simplicity,
 611 with degree d . Using Theorem 3.4, we have that the asymptotic covariance matrix
 612 equals:

$$613 \quad \mathbf{S} = a^2 \int_0^\infty e^{\Sigma v} \mathbf{S}_0 e^{\Sigma v} dv.$$

614 Here, $\mathbf{S}_0 = \left(h^2 \sigma_{\text{obs}}^2 + \frac{b^2}{a^2} d \sigma^2 \right) \mathbf{I}$; also, recall $\sigma^2 = \int_{-\infty}^\infty |\Psi(w)|^2 f(w) dw$, the effective
 616 communication noise per link. We assume that $f(w)$ has a zero mean and variance
 617 $\sigma_{\text{comm}}^2 = \int_{-\infty}^\infty w^2 f(w) dw$ that is finite. Also,

$$618 \quad \Sigma = \frac{1}{2} \mathbf{I} - a \left(h^2 \mathbf{I} + \frac{b}{a} \varphi'(0) \mathbf{L} \right),$$

619 where φ is given in (3.2). For the nonlinearity considered here, we have that

$$621 \quad \sigma^2 = 2 \int_0^{+B} w^2 f(w) dw + B^2 \left(1 - 2 \int_0^{+B} f(w) dw \right),$$

$$622 \quad \varphi'(0) = 2 \int_0^{+B} f(w) dw.$$

623 Denote by $\sigma_B^2 = \frac{1}{N} \text{Tr}(\mathbf{S})$ the average per-agent asymptotic variance. Analogously to
 625 (76)-(86) in [17], for $a > \frac{1}{2h^2}$ we get:

$$626 \quad \sigma_B^2 = \frac{a^2 h^2 \sigma_{\text{obs}}^2 + b^2 d \sigma^2}{N(2ah^2 - 1)} + \frac{a^2 h^2 \sigma_{\text{obs}}^2 + b^2 d \sigma^2}{N} \sum_{i=2}^N \frac{1}{2b\lambda_i \varphi'(0) + (2ah^2 - 1)}$$

628 We next analyze the values of σ_B^2 as $B \rightarrow 0$ and $B \rightarrow +\infty$.

629 For $B \rightarrow 0$, we have that $\sigma^2 \rightarrow 0$, $\varphi'(0) \rightarrow 0$ and

$$630 \quad \sigma_B^2 \rightarrow \frac{a^2 h^2 \sigma_{\text{obs}}^2}{2ah^2 - 1} =: \sigma_0^2.$$

632 That is, when $B \rightarrow 0$, we effectively have the case that each agent is working in
 633 isolation, hence not seeing the effect of the communication noise.

634 For $B \rightarrow +\infty$, we have that $\varphi'(0) \rightarrow 1$, $\sigma^2 \rightarrow \sigma_{\text{comm}}^2$ and

$$635 \quad \sigma_B^2 \rightarrow \frac{a^2 h^2 \sigma_{\text{obs}}^2 + b^2 d \sigma_{\text{comm}}^2}{N(2ah^2 - 1)} + \frac{a^2 h^2 \sigma_{\text{obs}}^2 + b^2 d \sigma_{\text{comm}}^2}{N} \sum_{i=2}^N \frac{1}{2b\lambda_i + (2ah^2 - 1)} =: \sigma_\infty^2.$$

637 This is the asymptotic variance of the linear \mathcal{LU} scheme in [17]. Note that, for any
 638 set of values of system parameters and any $a > \frac{1}{2h^2}$ and $b > 0$, there holds that

$$639 \quad (4.2) \quad \sigma_\infty^2 > \sigma_0^2$$

641 for a sufficiently large σ_{comm}^2 .

642 Assume from now on that (4.2) holds. It can be shown that there exists an optimal
 643 B , i.e., there exists B^* such that $B^* \in (0, +\infty)$ and $\inf_{B \in (0, +\infty)} \sigma_B^2 = \sigma_{B^*}^2$ (see Supple-

644 mentary material D).

645 Note that the above analysis generalizes also to the case when

$$646 \quad \sigma_{\text{comm}}^2 = \int_{-\infty}^{+\infty} w^2 f(w) dw = +\infty,$$

647

648 i.e., when the noise variance is $+\infty$. In this case, we have that $\sigma_\infty^2 = +\infty$, for the
 649 linear scheme and $\sigma_0^2 = \frac{a^2 h^2 \sigma_{\text{obs}}^2}{2ah^2 - 1}$ for the isolation scheme. It can be shown that
 650 $\inf_{B \in (0, +\infty)} \sigma_B^2$ is achieved at some $B^* \in (0, +\infty)$ (see Supplementary material D).

651 In order to demonstrate the results above, we minimize σ_B^2 and calculate B^* for a
 652 specific numerical example (see Figure 1a). We consider a sensor (agents) network
 653 with $N = 8$ agents, where the underlying topology is given by a regular graph with
 654 degree $d = 3$. We set innovation and consensus constants as $a = b = 1$, the observation
 655 parameter $h = 1$, and the true parameter $\theta^* = 1$. The observation noise for each
 656 sensor's measurements is standard normal, and the communication noise for each
 657 communication link has the following pdf

$$(4.3) \quad f(w) = \frac{\beta - 1}{2(1 + |w|)^\beta},$$

658 with $\beta = 2.05$. (This pdf's distribution has the infinite variance.) Figure 1b shows
 659 performance of the nonlinear consensus+innovations estimator (2.3) in terms of the
 660 estimated per-sensor mean squared error (MSE) across iterations, for the optimal B^*
 661 and for some sub-optimal choices of B , obtained through a Monte Carlo simulation.
 662 We can see that the scheme with B^* performs better than for the considered sub-
 663 optimal choices of B . Figure 1c shows that Monte Carlo estimate of the per-agent
 664 asymptotic variance, i.e., $\hat{S} = \frac{1}{N} \|\mathbf{x}^t - \mathbf{1}_N \otimes \theta^*\|^2 t$ matches well the corresponding
 665 theoretical value as per Theorem 3.4.
 667

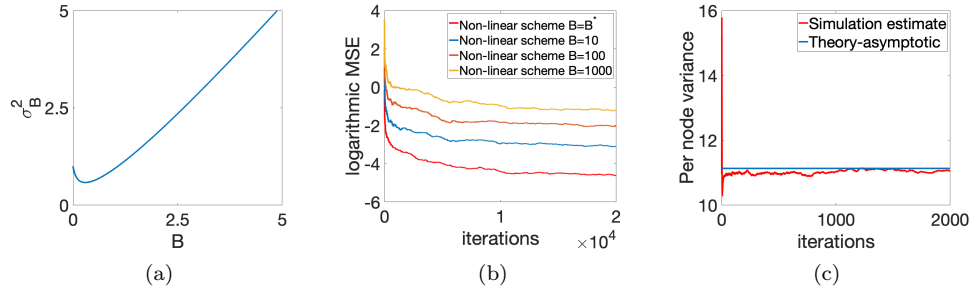


FIG. 1. (a) Per-agent asymptotic variance σ_B^2 versus B for the nonlinear consensus+ innovations estimator and the NL2 nonlinearity. (b) Monte Carlo-estimated per-sensor MSE error on logarithmic scale for the nonlinear consensus+innovations estimator with the NL2 nonlinearity for different choices of B . (c) Monte Carlo estimate of the per-agent asymptotic variance, and the corresponding theoretical value as per Theorem 3.4.

668 **Example 2:** We consider the same network and sensing models as in Example 1 and
 669 the heavy-tail communication noise distribution in (4.3). Furthermore, we assume
 670 that $\Psi(w) = \text{sign}(w)$ (the NL3 nonlinearity). For the \mathcal{LU} scheme, it can be shown
 671 that (see Supplementary material E):

$$672 \quad \sigma^2 = \sigma_{\text{comm}}^2 = \frac{2}{(\beta - 3)(\beta - 2)},$$

$$673 \quad \varphi'(0) = 1.$$

674 It can be shown here that the average per-agent asymptotic variance $\sigma_L^2 = \frac{1}{N} \text{Tr}(\mathbf{S})$
 675 for the \mathcal{LU} scheme is equal to
 676

$$(4.4) \quad \sigma_L^2 = \begin{cases} \infty & , \quad 2 < \beta \leq 3, \\ \frac{a^2 h^2 \sigma_{\text{obs}}^2 + b^2 d \sigma^2}{N(2ah^2 - 1)} + \frac{a^2 h^2 \sigma_{\text{obs}}^2 + b^2 d \sigma^2}{N} \sum_{i=2}^N \frac{1}{2b\lambda_i + (2ah^2 - 1)} & , \quad \beta > 3. \end{cases}$$

For $\beta > 3$, quantity σ_L^2 can be written as

$$(4.5) \quad \sigma_L^2 = A_L + B_L \frac{1}{(\beta - 3)(\beta - 2)},$$

where

$$A_L = \frac{a^2 h^2 \sigma_{\text{obs}}^2}{N(2ah^2 - 1)} + \frac{a^2 h^2 \sigma_{\text{obs}}^2}{N} \sum_{i=2}^N \frac{1}{2b\lambda_i + (2ah^2 - 1)},$$

$$B_L = 2 \left(\frac{b^2 d}{N(2ah^2 - 1)} + \frac{b^2 d}{N} \sum_{i=2}^N \frac{1}{2b\lambda_i + (2ah^2 - 1)} \right).$$

We next consider the nonlinear consensus+innovations scheme with the nonlinearity $\Psi(w) = \text{sign } w$. We have that

$$\sigma^2 = 1,$$

$$\varphi(a) = 2 \int_0^a f(w) dw,$$

which means that $\varphi'(a) = 2f(a)$ and $\varphi'(0) = 2f(0) = (\beta - 1)$. Hence, we have that the average per-agent asymptotic variance for the nonlinear scheme $\sigma_{\text{NL}}^2 = \frac{1}{N} \text{Tr}(S)$ is given by:

$$(4.6) \quad \sigma_{\text{NL}}^2 = \frac{a^2 h^2 \sigma_{\text{obs}}^2 + b^2 d \sigma^2}{N(2ah^2 - 1)} + \frac{a^2 h^2 \sigma_{\text{obs}}^2 + b^2 d \sigma^2}{N} \sum_{i=2}^N \frac{1}{4b\lambda_i f(0) + (2ah^2 - 1)},$$

which can be written in the form

$$(4.7) \quad \sigma_{\text{NL}}^2 = A_{\text{NL}} + B_{\text{NL}} \frac{P_{N-2}(\beta)}{\prod_{i=2}^N (\beta - \beta_i)},$$

where

$$A_{\text{NL}} = \frac{a^2 h^2 \sigma_{\text{obs}}^2 + b^2 d}{N(2ah^2 - 1)},$$

$$B_{\text{NL}} = \frac{a^2 h^2 \sigma_{\text{obs}}^2 + b^2 d}{N \prod_{i=2}^N 2b\lambda_i},$$

$$P_{N-2}(\beta) = \sum_{i=2}^N \prod_{\substack{j=2 \\ j \neq i}}^N 2b\lambda_j (\beta - \beta_j).$$

$$\beta_i = 1 - \frac{2ah^2 - 1}{2b\lambda_i}, \quad i = 2, \dots, N.$$

We next compare the average per-agent asymptotic variances for the linear consensus+innovations scheme and the nonlinear consensus+innovations scheme. From (4.4) it is obvious that $\sigma_{\text{NL}}^2 < \sigma_L^2$ for $\beta \in (2, 3]$. For $\beta > 3$, if $A_L \gg A_{\text{NL}}$ (see Supplementary material E), the linear scheme is worse than the nonlinear scheme for all $\beta > 3$. It is obvious that σ_L^2 decreases on interval $(3, \infty)$ and σ_{NL}^2 decreases on the interval (β_m, ∞) , where $\beta_m = \max_{i=2, \dots, N} \beta_i < 1$ is closest β_i to 1. Function $\sigma_L^2 = \sigma_L^2(\beta)$ has an asymptote at $\beta = 3$, and function $\sigma_{\text{NL}}^2 = \sigma_{\text{NL}}^2(\beta)$ at $\beta = \beta_m$, where $\beta_m < 3$, also,

712 A_L and A_{NL} are horizontal asymptotes for σ_L^2 and σ_{NL}^2 , respectively. Therefore, if
 713 A_L is much larger than A_{NL} , σ_L^2 is above σ_{NL}^2 for all $\beta > 3$. Moreover, if $A_L < A_{NL}$
 714 there exists $\beta^* > 3$ such that the average per-agent asymptotic variance is still better
 715 for the nonlinear than for the linear scheme for $\beta \in (2, \beta^*]$. Defining $k = \frac{\sigma_L^2}{\sigma_{NL}^2}$, it is
 716 possible to show that $k \rightarrow \infty$ as $\beta \rightarrow 3$, and $k \rightarrow \frac{A_L}{A_{NL}}$ as $\beta \rightarrow \infty$. Therefore, if
 717 $A_L < A_{NL}$, there exists β^* such that $\sigma_{NL}^2 < \sigma_L^2$ for all $\beta \in (2, \beta^*)$. In other words,
 718 there exists a threshold value $\beta^* > 3$, such that the nonlinear scheme outperforms the
 719 linear scheme for the “heavy-tail regime” $\beta \in (2, \beta^*)$, and the linear scheme performs
 720 better for $\beta > \beta^*$. To summarize, in Example 2, depending on sensing and network
 721 parameters, it holds that either the nonlinear scheme outperforms the linear one for
 722 all β , or there exists a threshold value β^* such that the nonlinear scheme is better
 723 than the linear one for $\beta \in (2, \beta^*)$. Figure 2 shows the ratio $k = \frac{\sigma_L^2}{\sigma_{NL}^2}$ versus β for
 724 the same sensing and network parameters as in Example 1. As it can be seen, there
 725 exists a threshold β^* , that here approximately equals $\beta^* = 3.9$, such that $k > 1$ for
 726 $\beta \in (2, \beta^*)$. On the other hand, for $\beta > \beta^*$, the ratio becomes smaller than one, which
 727 means that for the given numerical parameters, the linear scheme performs better for
 728 $\beta > \beta^*$. This is in accordance with the analysis that we provided above.

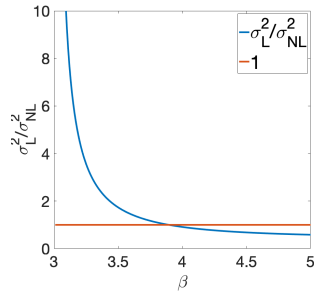


FIG. 2. Ratio $k = \frac{\sigma_L^2}{\sigma_{NL}^2}$ versus β for Example 2.

729 **4.2. Simulation examples.** In this section, we illustrate the performance of
 730 the proposed nonlinear consensus+innovations estimator for two different choices of
 731 the non-linear operator Ψ . For both nonlinearity choices, our method is compared
 732 with the corresponding linear consensus+innovations estimator \mathcal{LU} in [17], when the
 733 communication noise has probability distribution function given by (4.3).
 734 We consider a sensor network with $N = 40$ agents. The underlying topology is an
 735 instance of a random geometric graph. We use the same initialization $\mathbf{x}^0 = \mathbf{0}$ and same
 736 step sizes $\alpha_t = \frac{1}{t+1}$, $a = 1, b = 1$, for both the linear and the nonlinear estimators.
 737 Also, we assume that the observation noise is normally distributed, i.e., $n_i^t \sim \mathcal{N}(0, 1)$,
 738 for each t , for each i . The true parameter $\theta^* \in \mathbb{R}^{10}$ is generated randomly, where
 739 the entries of θ^* are drawn mutually independently from the uniform distribution
 740 on $[-10, 10]$. The observation vectors $\mathbf{h}_i \in \mathbb{R}^{10}$ are also generated at random, for
 741 which the condition 4 of Assumption 2.3 is true. We use the communication noise
 742 pdf in (4.3) with $\beta = 2.05$. Note that, in this case, the communication noise has an
 743 infinite variance.

744 Figure 4 compares the linear \mathcal{LU} estimator in [17] with the nonlinear estima-
 745 tor (2.3) with $\Psi(w)$ given in (4.1) for $B = 5$. Figure 3 shows the comparison between
 746 \mathcal{LU} and [17] with $\Psi(w) = \text{sign}(w)$. Both Figures show the iteration counter t at the

747 x -axis and a Monte-Carlo estimate of the average mean square error (MSE) across
 748 agents on the y -axis. We can see that, as predicted by our theory, the nonlinear
 749 estimator, for both nonlinearity choices, persistently decreases MSE along iterations,
 750 despite the fact that the communication noise has an infinite variance. At the same
 751 time, \mathcal{LU} fails to produce a useful estimation result.

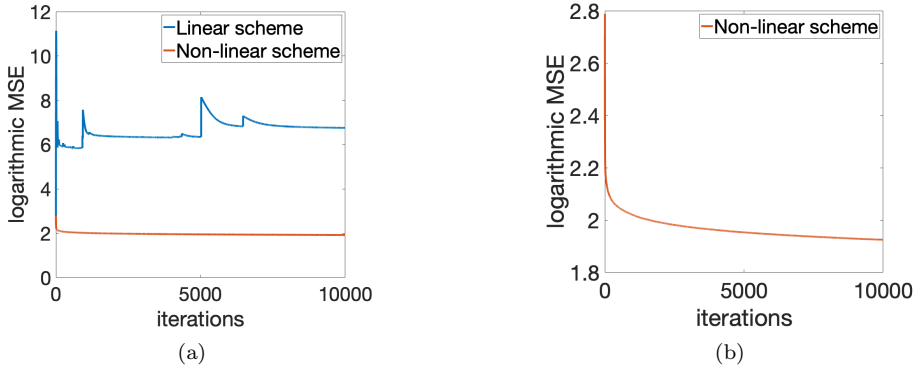


FIG. 3. Monte-Carlo average per-agent MSE estimate versus iteration counter on logarithmic scale for the proposed nonlinear estimator (2.3) with the nonlinearity in (4.1) for $B = 5$ and the linear \mathcal{LU} scheme in [17].

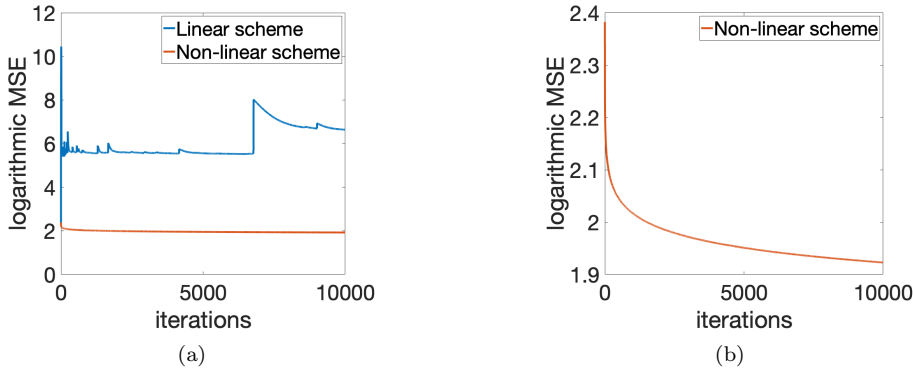


FIG. 4. Monte-Carlo average per-agent MSE estimate versus iteration counter on logarithmic scale for the proposed nonlinear estimator (2.3) with the nonlinearity $\Psi(w) = \text{sign}(w)$ and the linear \mathcal{LU} scheme in [17].

752 **5. Conclusion.** We studied consensus+innovations distributed estimation in the
 753 presence of impulsive, heavy-tail communication noise. To combat the impulsive
 754 communication noise, we introduce for the first time a general nonlinearity in the
 755 consensus update for consensus+innovations distributed estimation. We establish al-
 756 most sure convergence of the nonlinear consensus+innovations estimator to the true
 757 parameter, prove its asymptotic normality, and explicitly evaluate the corresponding
 758 asymptotic variance. We compare the proposed nonlinear estimator with conventional
 759 consensus+innovation estimators that utilize linear consensus update. Analytical and

760 numerical examples demonstrate significant gains of introducing consensus nonlinear-
 761 ity in low SNR (high communication noise) regimes. Most notably, we demonstrate
 762 that, when the communication noise has infinite variance, the proposed nonlinear con-
 763 sensus+innovations estimator is strongly consistent (converges almost surely), while
 764 the corresponding linear counterpart provides a sequence of estimators with infinite
 765 variance.

766 **Acknowledgments** The work of D. Bajovic, D. Jakovetic and M. Vukovic is
 767 supported by the Ministry of Education, Science and Technological Development,
 768 Republic of Serbia. Moreover, the work of D. Bajovic and D. Jakovetic is supported
 769 by the European Union’s Horizon 2020 Research and Innovation program under grant
 770 agreement No. 957337. The paper reflects only the view of the authors and the
 771 Commission is not responsible for any use that may be made of the information it
 772 contains.

773

REFERENCES

- 774 [1] S. AL-SAYED, A. M. ZOUBIR, AND A. H. SAYED, *Robust distributed estimation by networked*
 775 *agents*, IEEE Transactions on Signal Processing, 65 (2017), pp. 3909–3921.
- 776 [2] D. BAJOVIC, D. JAKOVETIC, J. M. MOURA, J. XAVIER, AND B. SINOPOLI, *Large deviations per-*
 777 *formance of consensus+ innovations distributed detection with non-gaussian observations*,
 778 IEEE Transactions on Signal Processing, 60 (2012), pp. 5987–6002.
- 779 [3] D. BAJOVIC, D. JAKOVETIC, J. XAVIER, B. SINOPOLI, AND J. MOURA, *Distributed detection*
 780 *via gaussian running consensus: Large deviations asymptotic analysis*, Signal Processing,
 781 IEEE Transactions on, 59 (2011), pp. 4381 – 4396.
- 782 [4] W. BEN-AMEUR, P. BIANCHI, AND J. JAKUBOWICZ, *Robust distributed consensus using total*
 783 *variation*, IEEE Transactions on Automatic Control, 61 (2016), pp. 1550–1564.
- 784 [5] M. CAO, A. MORSE, AND B. ANDERSON, *Reaching a consensus in a dynamically changing en-*
 785 *vironment: A graphical approach*, SIAM J. Control and Optimization, 47 (2008), pp. 575–
 786 600.
- 787 [6] Y. CHEN, S. KAR, AND J. MOURA, *Resilient distributed field estimation*, SIAM Journal on
 788 Control and Optimization, 58 (2020), pp. 1429–1456.
- 789 [7] S. CHOUVARDAS, K. SLAVAKIS, AND S. THEODORIDIS, *Adaptive robust distributed learning in*
 790 *diffusion sensor networks*, Signal Processing, IEEE Transactions on, 59 (2011), pp. 4692 –
 791 4707.
- 792 [8] L. CLAVIER, T. PEDERSEN, I. LARRAD, M. LAURIDSEN, AND M. EGAN, *Experimental evidence*
 793 *for heavy tailed interference in the IoT*, IEEE Communications Letters, 25 (2021), pp. 692–
 794 695.
- 795 [9] S. DASARATHAN, C. TEPEDELENLIOĞLU, M. K. BANAVAR, AND A. SPANIAS, *Robust consensus*
 796 *in the presence of impulsive channel noise*, IEEE Transactions on Signal Processing, 63
 797 (2015), pp. 2118–2129.
- 798 [10] F. FAGNANI AND S. ZAMPIERI, *Average consensus with packet drop communication*, in Proceed-
 799 ings of the 45th IEEE Conference on Decision and Control, 2006, pp. 1007–1012.
- 800 [11] M. HUANG AND J. MANTON, *Coordination and consensus of networked agents with noisy mea-*
 801 *surements: Stochastic algorithms and asymptotic behavior*, SIAM J. Control and Opti-
 802 mization, 48 (2009), pp. 134–161.
- 803 [12] B. HUGHES, *Alpha-stable models of multiuser interference*, in 2000 IEEE International Sympo-
 804 sium on Information Theory (Cat. No.00CH37060), 2000, pp. 383–.
- 805 [13] J. ILOW AND D. HATZINAKOS, *Analytic alpha-stable noise modeling in a poisson field of in-*
 806 *terferers or scatterers*, Signal Processing, IEEE Transactions on, 46 (1998), pp. 1601 –
 807 1611.
- 808 [14] D. JAKOVETIC, J. M. F. MOURA, AND J. XAVIER, *Distributed detection over noisy networks:*
 809 *Large deviations analysis*, IEEE Transactions on Signal Processing, 60 (2012), pp. 4306–
 810 4320.
- 811 [15] S. KAR AND J. MOURA, *Asymptotically efficient distributed estimation with exponential family*
 812 *statistics*, IEEE Transactions on Information Theory, 60 (2014), pp. 4811–4831.
- 813 [16] S. KAR, J. MOURA, AND H. V. POOR, *Distributed linear parameter estimation: Asymptoti-*
 814 *cally efficient adaptive strategies*, SIAM Journal on Control and Optimization, 51 (2013),
 815 pp. 2200–2229.

- 816 [17] S. KAR, J. M. F. MOURA, AND K. RAMANAN, *Distributed parameter estimation in sensor*
817 *networks: Nonlinear observation models and imperfect communication*, IEEE Transactions
818 on Information Theory, 58 (2012), pp. 3575–3605.
- 819 [18] U. A. KHAN, S. KAR, AND J. M. F. MOURA, *Distributed average consensus: Beyond the realm*
820 *of linearity*, in 2009 Conference Record of the Forty-Third Asilomar Conference on Signals,
821 Systems and Computers, 2009, pp. 1337–1342.
- 822 [19] S. KUMAR, U. K. SAHOO, A. K. SAHOO, AND D. P. ACHARYA, *Diffusion minimum-wilcoxon-*
823 *norm over distributed adaptive networks: Formulation and performance analysis*, Digital
824 Signal Processing, 51 (2016), pp. 156–169.
- 825 [20] A. LALITHA, T. JAVIDI, AND A. D. SARWATE, *Social learning and distributed hypothesis testing*,
826 IEEE Transactions on Information Theory, 64 (2018), pp. 6161–6179.
- 827 [21] Z. LI AND S. GUAN, *Diffusion normalized huber adaptive filtering algorithm*, Journal of the
828 Franklin Institute, 355 (2018), pp. 3812–3825.
- 829 [22] Q. LIU AND A. IHLER, *Distributed estimation, information loss and exponential families*, 2014.
- 830 [23] C. LOPES AND A. SAYED, *Diffusion least-mean squares over adaptive networks: Formulation*
831 *and performance analysis*, IEEE Transactions on Signal Processing, 56 (2008), pp. 3122–
832 3136.
- 833 [24] G. MATEOS, I. SCHIZAS, AND G. GIANNAKIS, *Distributed recursive least-squares for consensus-*
834 *based in-network adaptive estimation*, Signal Processing, IEEE Transactions on, 57 (2009),
835 pp. 4583 – 4588.
- 836 [25] V. MATTA, P. BRACA, S. MARANO, AND A. H. SAYED, *Diffusion-based adaptive distributed*
837 *detection: Steady-state performance in the slow adaptation regime*, IEEE Transactions on
838 Information Theory, 62 (2016), pp. 4710–4732.
- 839 [26] S. MODALAVALASA, U. SAHOO, A. SAHOO, AND S. BARAHA, *A review of robust distributed esti-*
840 *mation strategies over wireless sensor networks*, Signal Processing, 188 (2021), p. 108150.
- 841 [27] A. NEDIC, A. OLSHEVSKY, AND C. A. URIBE, *Nonasymptotic convergence rates for cooperative*
842 *learning over time-varying directed graphs*, in 2015 American Control Conference (ACC),
843 IEEE, 2015, pp. 5884–5889.
- 844 [28] M. B. NEVEL’SON AND R. Z. HAS’ MINSKII, *Stochastic approximation and recursive estimation*,
845 vol. 47, American Mathematical Soc., 1976.
- 846 [29] B. POLYAK AND Y. TSYPKIN, *Adaptive estimation algorithms: Convergence, optimality, stabil-*
847 *ity*, Automation and Remote Control, 1979 (1979).
- 848 [30] A. PRASAD, A. S. SUGGALA, S. BALAKRISHNAN, AND P. RAVIKUMAR, *Robust estimation via*
849 *robust gradient estimation*, Journal of the Royal Statistical Society: Series B (Statistical
850 Methodology), 82 (2020), pp. 601–627.
- 851 [31] S. RAM, V. VEERAVALLI, AND A. NEDIC, *Distributed and Recursive Parameter Estimation*,
852 Springer Science & Business Media, 2009, pp. 17–38.
- 853 [32] B. SELIM, M. S. ALAM, V. CARVALHO, G. KADDOUM, AND B. L. AGBA, *Noma-based iot net-*
854 *works: Impulsive noise effects and mitigation*, IEEE Communications Magazine, 58 (2020),
855 pp. 69–75.
- 856 [33] S. STANKOVIC, M. BEKO, AND M. STANKOVIC, *A robust consensus seeking algorithm*, in IEEE
857 EUROCON 2019-18th International Conference on Smart Technologies, 2019, pp. 1–6.
- 858 [34] S. SUNDARAM AND B. GHARESIFARD, *Consensus-based distributed optimization with malicious*
859 *nodes*, in 2015 53rd Annual Allerton Conference on Communication, Control, and Com-
860 puting (Allerton), 2015, pp. 244–249.
- 861 [35] S. THEODORIDIS, K. SLAVAKIS, AND I. YAMADA, *Adaptive learning in a world of projections*,
862 Signal Processing Magazine, IEEE, 28 (2011), pp. 97 – 123.
- 863 [36] F. WEN, *Diffusion least mean p -power algorithms for distributed estimation in alpha-stable*
864 *noise environments*, Electronics Letters, 49 (2013).
- 865 [37] X. YANG AND A. PETROPULU, *Co-channel interference modeling and analysis in a poisson*
866 *field of interferers in wireless communications*, IEEE Transactions on Signal Processing,
867 51 (2003), pp. 64–76.
- 868 [38] X. ZHAO, S.-Y. TU, AND A. H. SAYED, *Diffusion adaptation over networks under imperfect*
869 *information exchange and non-stationary data*, IEEE Transactions on Signal Processing,
870 60 (2012), pp. 3460–3475.