

Distributed Adaptive Spatial Filtering with Inexact Local Solvers

Charles Hovine , Alexander Bertrand 

Abstract—The Distributed Adaptive Signal Fusion (DASF) framework is a meta-algorithm for computing data-driven spatial filters in a distributed sensing platform with limited bandwidth and computational resources, such as a wireless sensor network. The convergence and optimality of the DASF algorithm has been extensively studied under the assumption that an exact, but possibly impractical solver for the local optimization problem at each updating node is available. In this work, we provide convergence and optimality results for the DASF framework when used with an inexact, finite-time solver such as (proximal) gradient descent or Newton’s method. We provide sufficient conditions that the solver should satisfy in order to guarantee convergence of the resulting algorithm, and a lower bound for the convergence rate. We also provide numerical simulations to validate these theoretical results.

Index Terms—Distributed signal processing, Wireless sensor networks, Adaptive spatial filtering, Optimization.

I. INTRODUCTION

A wireless sensor network (WSN) is a networked collection of sensor nodes equipped with communication and computing capabilities. Each node typically senses a single or multi-channel signal, which can either be transmitted to a fusion center for processing, or be processed *in-network*. Spatial filtering is a common processing task in such networks, allowing the extraction of some target signal from the aggregated sensor signals. A spatial filter typically consists in a linear combination of the sensor signals, producing a lower-dimensional output signal which is optimal in some sense, for example, in terms of signal-to-noise ratio, output power, or correlation with a target signal. Common examples include principal component analysis (PCA) [1], Wiener filtering [2], canonical correlation analysis (CCA) [3], linearly constrained minimum variance beamforming (LCMV) [4] and Max-SNR filtering [5].

The most straightforward way to compute such filters is to aggregate the sensor data at a fusion center, where an off-the-shelf algorithm can be used to process the data. This convenience comes with several drawbacks. Firstly, the transmission of the sensor nodes’ raw sensor signals are likely to

become the main energy bottleneck of the system [6], in particular in networks where multi-hop data relaying is required, resulting in shorter battery life, making remote deployments impractical. This poses a particular challenge for applications that generate continuous streams of high-throughput sensor data, such as audio in acoustic sensor networks [7], video in CCTV systems or video sensor networks [8], biomedical signals in EEG¹ sensor networks [9], [10], or radio signals in multistatic radars [11], [12]. Secondly, the fusion center needs to have sufficient computing capabilities to handle the continuous processing of the stream of high-dimensional input signals. This comes with scalability issues, as the increase of the number of nodes or input signals will eventually lead to unrealistic hardware requirements at the fusion center in terms of bandwidth, memory, and compute. Finally, the use of a fusion center introduces a single point of failure [13], which can be prohibitive if maintenance actions are difficult or costly, such as in the case of remote deployments. The drawbacks associated with centrally processing the data motivate the use of distributed signal processing algorithms that distribute the computational burden amongst the nodes and favor local processing over data transmission.

We can identify two classes of distributed signal processing algorithms, associated with two corresponding kinds of distributed datasets. The first kind covers algorithms operating on the observations of a stochastic signal $\mathbf{y}(t) \in \mathbb{R}^M$, whose observations are distributed across the nodes, such that each node has access to all M entries of $\mathbf{y}(t)$, but only a subset of the realizations of the random data associated with those entries. This allows each node to locally estimate the covariance structure of the full signal. The optimality criterion can in this case typically be expressed as a sum of local per-node objectives, and the distributed algorithm will usually consist in performing some local processing on the node, and then exchanging some low-dimensional vector containing intermediate estimates of the optimization variables (rather than signal samples). Typical examples of algorithms suited for such datasets are the alternating direction method of multipliers (ADMM) [14], consensus and diffusion strategies [15], [16], and many of the techniques studied by the federated learning community [17].

The second kind of distributed algorithms deals with datasets where the entries (or channels) of $\mathbf{y}(t)$ are spread across the nodes, such that no individual node can estimate the correlation structure of the network-wide signal $\mathbf{y}(t)$ (i.e. the signal consisting in the concatenation of all the channels of the per-node signals). Such algorithms typically

This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No. 802895 and No. 101138304) and from the Flemish Government under the “Onderzoeksprogramma Artificiële Intelligentie (AI) Vlaanderen” programme. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or ERC. Neither the European Union nor the ERC can be held responsible for them.

Charles Hovine and Alexander Bertrand are with the STADIUS Center for Dynamical Systems, Signal Processing and Data Analytics and with the Leuven.AI institute for Artificial Intelligence at KU Leuven, Leuven 3001, Belgium (e-mails: {charles.hovine, alexander.bertrand}@esat.kuleuven.be).

¹Electro-encephalography.

require that signal samples (rather than low-dimensional parameter vectors) are shared between the nodes, in order to learn the inter-node statistics and to properly evaluate the optimality criterion. Distributed spatial filtering problems or regression/classification problems with distributed features fit in this second class, and are more difficult to solve if these inter-node statistics are not known a-priori. For streaming data or data sets where the number of observations of $\mathbf{y}(t)$ is much larger than the number of channels M , the typical work-horse distributed methods (ADMM, consensus, diffusion, etc.) do not straightforwardly apply [18], or result in highly inefficient multi-rate communication schemes with nested iterations [19], [20]. Therefore, they are typically solved with ad-hoc problem-specific algorithms, due to the absence of a one-fit-all off-the-shelf solution [4], [20]–[26]. [27] describes an algorithm to solve such distributed-features problem, but requires a central node (although most of the computational load is distributed amongst the nodes). In addition, it is limited to smooth and unconstrained problems.

The Distributed Adaptive Signal Fusion (DASF) framework [28]–[31] aims to bridge the gap by providing a scalable and bandwidth-efficient algorithm that adaptively computes the outputs of data-driven spatial filters satisfying some optimality criterion, in a WSN or other distributed setting with constrained resources. Given an existing solver for the particular optimization problem associated with the filter of interest, the DASF algorithm can adaptively compute the desired filter along with the filtered signal by relying on the exchange of low-dimensional compressed samples between the nodes, and the repeated computation of the solution of a lower-dimensional version of the original centralized optimization problem at each node. The DASF framework is mostly problem-agnostic, and most traditional linear spatial filtering or estimation problems can be solved by plugging-in the appropriate (and unmodified) solver for the corresponding centralized problem. However, the original DASF framework was shown to converge under the assumption that the solver produces an (near-)exact solution every time it is called. In practice though, an exact solver might not be available, either due to the nature of the problem (e.g. non-convexity), or simply because obtaining an exact solution would be too expensive or time consuming, hindering the adaptivity of the algorithm. In addition, the convergence of DASF relies on hard-to-check assumptions regarding the parametric continuity of the problem. In this work, we show that an exact solution is not required for optimality, and that a few iterations of an iterative solver are in practice sufficient to ensure convergence to an “interesting” filter (i.e. satisfying some relaxed optimality condition such as mere stationarity with regards to the optimality criterion). We are also able to eliminate some of the original convergence assumptions, making the algorithm applicable to a broader set of problems.

The paper is organized as follows. Section II describes the general objective, as well as the family of problems covered by DASF. Section III gives a brief overview of the original DASF algorithm. Section IV describes how DASF can be modified to work with an inexact iterative solver, rather than an exact one. This section describes our main contribution, including a proof

of the convergence of DASF used with such a solver. Section V offers a validation of our theoretical results via numerical simulations. We conclude with a brief discussion in Section VI.

II. PROBLEM STATEMENT

We consider a WSN consisting of K nodes, each sensing an M_k -dimensional stochastic signal $\mathbf{y}_k(t)$, where t denotes a sample index (usually related to a time index). We denote the M -dimensional network-wide sensor signal $\mathbf{y}(t) = [\mathbf{y}_1^T(t), \dots, \mathbf{y}_K^T(t)]^T$ with values in \mathbb{R}^M , with $M = \sum_k M_k$. We assume that \mathbf{y} is short-term stationary and ergodic, such that its statistics can be computed over short-term sample batches. The DASF framework assumes that the nodes in the network collaborate to compute an optimal linear spatial filter $X \in \mathbb{R}^{M \times Q}$ in a bandwidth-efficient manner. Similarly to \mathbf{y} , we define the per-node partition of X as $X = [X_1^T, \dots, X_K^T]^T$, where we refer to X_k as the “local” filter associated with node k . The Q -channel filtered signal $\mathbf{z}(t) \triangleq X^T \mathbf{y}(t)$, and hence the filter X should satisfy some optimality criterion, which can be expressed as [28]

$$\begin{aligned} X^* \in \operatorname{argmin}_{X \in \mathbb{R}^{M \times Q}} \varphi(X^T \mathbf{y}(t), X^T B) \\ \text{s.t. } \quad \forall j \in \mathcal{J}_I, \eta_j(X^T \mathbf{y}(t), X^T D_j) \leq 0, \\ \quad \forall j \in \mathcal{J}_E, \eta_j(X^T \mathbf{y}(t), X^T D_j) = 0 \end{aligned} \quad (1)$$

where the matrices B and D_j are deterministic matrices known by every node, φ is a smooth real-valued function encoding some design objective for the filter output, \mathcal{J}_I and \mathcal{J}_E are the sets of inequality and equality constraints indices, respectively, and the η_j are smooth functions enforcing some hard constraints on the filter outputs and/or filter coefficients.

In order for the above functions to be real-valued, we assume that they implicitly contain an operator turning the random argument $X^T \mathbf{y}(t)$ into a deterministic one, by applying, e.g., an expectation operator. For the purpose of mathematical tractability, we assume for the theoretical analysis of the algorithm that $\mathbf{y}(t)$ is stationary (i.e. that its statistics are independent of time) and ergodic (i.e. its statistics can be estimated using samples averages), and its statistics are assumed to be perfectly estimated at any point in time. However, the expectation operators will in practice be estimated by temporal averages of $\mathbf{y}(t)$. The impact of such an approximation on the algorithm’s performance is outside of the scope of this paper, we therefore refer the interested reader to the stochastic optimization literature [32] for details.

The peculiar structure of problem (1) ensures that X always appears either in an inner product with $\mathbf{y}(t)$ or in an inner product with some pre-determined matrices (B or D_j). Most traditional spatial filtering or linear estimation problems satisfy this structure (two examples are given hereafter, and more extensive illustrations can be found in [28]). The DASF algorithm will exploit this inner-product structure to achieve a bandwidth reduction. Note that the framework also allows for complex-valued signals, in which case all transpose operators should be replaced with conjugate (a.k.a. Hermitian) transpose operators. Finally, it is noted that we only include a single \mathbf{y}

and B -matrix in the loss function φ , and a single D -matrix in each constraint. This is for the sake of notational convenience, as the DASF framework also allows multiple instances of each (details in [28]).

In order to also cover the treatment of non-smooth objective functions, we also allow the optimality criterion to be expressed as [30], [31]

$$\begin{aligned} X^* \in \operatorname{argmin}_{X \in \mathbb{R}^{M \times Q}} & \varphi(X^T \mathbf{y}(t), X^T B) + \gamma(X^T A) \\ \text{s.t. } \forall k \in \mathcal{K}, & \\ \forall j \in \mathcal{J}_I^k, & \eta_j(X_k^T \mathbf{y}_k(t), X_k^T D_{j,k}) \leq 0, \\ \forall j \in \mathcal{J}_E^k, & \eta_j(X_k^T \mathbf{y}_k(t), X_k^T D_{j,k}) = 0 \end{aligned} \quad (2)$$

where the main differences with (1) are the restriction of the constraints to be per-node block-separable, i.e., the constraints index sets \mathcal{J}_I^k and \mathcal{J}_E^k describe the constraints for a specific node k , and the functions η_j associated with a particular node k can only depend on the local filter X_k associated with that node, and the addition of γ , a possibly non-smooth function. The function γ is also required to be per-node block-separable, i.e., there exist matrices A_k and functions γ_k such that

$$\gamma(X^T A) = \sum_{k \in \mathcal{K}} \gamma_k(X_k^T A_k), \quad (3)$$

which also implies that A in (3) must be a block-diagonal matrix with blocks A_k . This allows for most of the typical regularizers, including the ℓ_1 , ℓ_2 , and the $\ell_{1,2}$ norm. Note that for simplicity, we previously assumed in [31] that γ was a convex function, but this assumption is here relaxed, and γ is thus allowed to be non-convex.

We refer to the parametric optimization problem defined by (1) as $\mathbb{P}_S(\mathbf{y}(t), B, \mathcal{D})$, and $\mathbb{P}_{NS}(\mathbf{y}(t), B, \mathcal{D}, A)$ for (2). Here the set \mathcal{D} denotes the collection of all the matrices D_j or $D_{j,k}$. These generic problem formulations encompass a wide range of well-known spatial filtering problems, including PCA, CCA, LCMV, and many more [28]. For example, we can express a PCA filter as

$$\begin{aligned} \min_X & -\operatorname{Tr}(X^T \mathbb{E}\{\mathbf{y}\mathbf{y}^T\} X) \\ \text{s.t. } & X^T X = I, \end{aligned} \quad (4)$$

where $\mathbb{E}\{\cdot\}$ denotes the expectation operator. A D matrix hides in the constraint, which can equivalently be written

$$X^T D D^T X = I, \quad D = I. \quad (5)$$

This might seem superfluous, but D actually plays an important algorithmic role, as its presence will allow the local sub-problems solved during the application of DASF to keep the same structure $\mathbb{P}_S(\mathbf{y}(t), B, \{D\})$ as the original one, but where $D \neq I$. The $\ell_{1,2}$ -regularized multi-channel Wiener filter [33] is a non-smooth example that fits in the class of problems defined by (1)-(2):

$$\min_X \mathbb{E} \left\{ \|X^T \mathbf{y} - \mathbf{d}\|_F^2 \right\} + \sum_k \|X_k\|_F \quad (6)$$

where $\|\cdot\|_F$ denotes the Frobenius norm, and $\mathbf{d}(t)$ is a known multichannel signal taking values in \mathbb{R}^Q . The filter produces an estimate of the signal $\mathbf{d}(t)$ from the measurements $\mathbf{y}(t)$, and

the non-smooth regularization term encourages nodes that do not significantly contribute to the objective to set their filters to zero. Note that X does not appear in an inner product in the second term of (6), yet it fits (1)-(2) when explicitly writing the Frobenius norm as $\|X_k^T A_k\|_F$ with $A_k = I$.

III. DASF OVERVIEW

In this section, we briefly review the DASF algorithm. We refer to [28] for more details and illustrative examples.

The DASF algorithm collaboratively updates the estimate of the optimal filter X^* and tracks the filtered signal \mathbf{z} by solving a local “miniature” version of (1) or (2) at a particular so-called *updating node* whose role is taken by a different node at every iteration. This makes the DASF algorithm “plug-and-play”² as it only requires a solver for the centralized problem (1) or (2), which is directly re-used as a local “sub-solver” in the compressed local problems in the different iterations of the distributed algorithm [28].

As the focus of this paper is the local sub-solver used by DASF, which is largely independent of the rest of the algorithmic framework, we can afford to limit our description of the algorithm to fully-connected networks. It is important to note that this is purely for intelligibility and notational convenience, but without loss of generality. The results described in this paper can be straightforwardly extended to arbitrary topology networks, as in [28], [29] (for the smooth case) and [31] (for the non-smooth case) and the convergence proof in Section IV-C is kept sufficiently generic to also cover the case of arbitrary-topology networks. As a result, the extension of DASF to arbitrary topologies that is detailed in [28] and [31] straightforwardly applies to the inexact version of DASF introduced in this paper, yet it is omitted here to reduce overlap and avoid the extra layer of notational complexity that it would incur.

An iteration i of the DASF algorithm is fully defined by the current updating node index q^i (we often drop the iteration index when the iteration is clear from the context) and the current estimate of the filter X^i , and consists of the following three steps (after a random initialization of X^0):

1. Data Aggregation At the beginning of a new iteration i , a new updating node q is selected. Each node k collects a new batch of N samples of $\mathbf{y}(t)$ for $t = iN, \dots, (i+1)N - 1$, and compresses these into Q -dimensional³ samples according to

$$\widehat{\mathbf{y}}_k^i \triangleq X_k^{iT} \mathbf{y}_k. \quad (7)$$

²A Matlab and Python toolbox implementing this plug-and-play functionality is available in [34].

³Here we assume that $Q < M_k$, otherwise there is no compression possible at node k . Nodes for which $M_k \leq Q$ can add their channels \mathbf{y}_k to the channels of a neighboring node to create a virtual “super-node” within the DASF algorithm. Note that in multi-hop network topologies, a bandwidth reduction can even be achieved if $Q > M_k$ in all nodes (when compared to relayed data aggregation). We refer to [28] for more details.

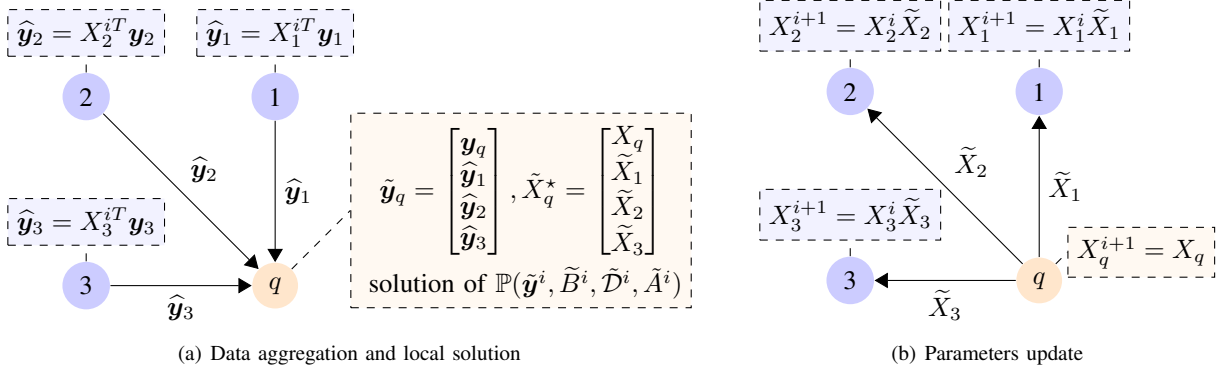


Fig. 1: Overview of a single iteration of the DASF Algorithm in fully-connected networks, where node q is the updating node. The matrices B_k , $D_{j,k}$ and A_k are omitted for readability, but their treatment is the same as \mathbf{y}_k .

The compressed N -samples batch is then transmitted to the updating node, along with the compressed matrices

$$\hat{D}_{j,k} \triangleq X_k^{iT} D_{j,k} \quad (8)$$

$$\hat{A}_k^i \triangleq X_k^{iT} A_k \quad (9)$$

$$\hat{B}_k^{iT} \triangleq X_k^{iT} B_k \quad (10)$$

(typically of negligible size compared to the transmission of (7) when the sample batch size N is large). Here, B_k is the block of B associated with the block X_k of X . In order to offer a unifying description of the algorithm for the smooth and non-smooth case, we also define the same relationship between $D_{j,k}$ and D_j (i.e. in the smooth case $D_{j,k}$ is the block of D_j associated with X_k in the product $X^T D_j$), and in order to be consistent with the definitions of problems (1) and (2), we require that $D_{j,k} = 0$ if $j \notin \mathcal{J}_E^k \cup \mathcal{J}_I^k$ in the non-smooth case (i.e. block-separable constraints), and $A_k = 0$ in the smooth-case (i.e. smooth objective function).

2. Local Solution The updating node constructs a *local* view $\tilde{\mathbf{y}}_q$ of \mathbf{y} by concatenating the received signals with its own sensor signals such that

$$\tilde{\mathbf{y}}^i = [\mathbf{y}_q^T \quad \hat{\mathbf{y}}_1^{iT} \quad \cdots \quad \hat{\mathbf{y}}_{q-1}^{iT} \quad \hat{\mathbf{y}}_{q+1}^{iT} \quad \cdots \quad \hat{\mathbf{y}}_K^{iT}]^T. \quad (11)$$

Similarly, it constructs

$$\begin{aligned} \tilde{A}^i &= \text{BlockDiag}(A_q, \hat{A}_1^i, \dots, \hat{A}_{q-1}^i, \hat{A}_{q+1}^i, \dots, \hat{A}_K^i), \\ \tilde{B}^i &= \begin{bmatrix} B_q^T & \hat{B}_1^{iT} & \cdots & \hat{B}_{q-1}^{iT} & \hat{B}_{q+1}^{iT} & \cdots & \hat{B}_K^{iT} \end{bmatrix}^T, \text{ and} \\ \tilde{D}_j^i &= \begin{bmatrix} D_{j,q}^T & \hat{D}_{j,1}^{iT} & \cdots & \hat{D}_{j,q-1}^{iT} & \hat{D}_{j,q+1}^{iT} & \cdots & \hat{D}_{j,K}^{iT} \end{bmatrix}^T \end{aligned} \quad (12)$$

The updating node then obtains the local solution \tilde{X}^* by solving $\mathbb{P}(\tilde{\mathbf{y}}^i, \tilde{B}^i, \tilde{D}^i, \tilde{A}^i)$ (which here refers to either \mathbb{P}_S or \mathbb{P}_{NS}) for the received N -sample batch, where, similarly to \mathcal{D} , $\tilde{\mathcal{D}}^i$ denotes the collection of \tilde{D}_j^i , which can be viewed as a low-dimensional instance of the original optimization problem. Note that the solver for the network-wide problem $\mathbb{P}(\mathbf{y}, B, \mathcal{D}, A)$ can thus also be used as-is to solve the ‘‘compressed’’ problem $\mathbb{P}(\tilde{\mathbf{y}}^i, \tilde{B}^i, \tilde{\mathcal{D}}^i, \tilde{A}^i)$ at updating node q [28].

3. Parameters Update The updating node partitions the local solution \tilde{X}^* as⁴

$$\tilde{X}^* = [X_q^{*T}, \tilde{X}_1^{*T}, \dots, \tilde{X}_{q-1}^{*T}, \tilde{X}_{q+1}^{*T}, \dots, \tilde{X}_K^{*T}]^T. \quad (13)$$

It then updates its own block of the filter as

$$X_q^{i+1} \leftarrow X_q^*, \quad (14)$$

and transmits each \tilde{X}_k to its corresponding node, which in turn updates its local filter as

$$X_k^{i+1} \leftarrow X_k^i \tilde{X}_k^*. \quad (15)$$

The three steps are illustrated in Figure 1 and the full procedure is described by Algorithm 1.

As a different batch of N samples is used at each iteration, the DASF algorithm produces the filtered signal

$$\mathbf{z}^{i+1} \triangleq X^{i+1T} \mathbf{y} = \sum_k X_k^{i+1T} \mathbf{y}_k = \tilde{X}^{*T} \tilde{\mathbf{y}}^i \quad (16)$$

for each N -samples block, while at the same time improving the estimate of the optimal spatial filter X^* , such that each new block of the filtered signal is closer to the desired filtered signal (under the stationarity assumption). In other words, DASF acts as a time-recursive block-adaptive filter that continually adapts itself over time to the (possibly changing) statistics of $\mathbf{y}(t)$ [28]. Note that the last equality implies that each updating node can locally produce the new estimate \mathbf{z}^{i+1} without additional data exchange.

Provided that the local solution is obtained using an exact solver, the DASF algorithm and its extension to arbitrary network topologies have been shown to converge to a stationary point of the original problem (under some technical conditions, details omitted) [29], [31].

IV. DASF WITH INEXACT LOCAL SOLVERS

The algorithm described in the previous section might in practice be difficult to implement, as the computational cost to find a global minimizer of $\mathbb{P}(\tilde{\mathbf{y}}^i, \tilde{B}^i, \tilde{\mathcal{D}}^i, \tilde{A}^i)$ at the updating node might be prohibitive, requiring either too much

⁴Note that $X_q^* \in \mathbb{R}^{M_q \times Q}$ and $\tilde{X}_k^* \in \mathbb{R}^{Q \times Q}$ for all other $k \neq q$. This follows from the dimensions of the partitioning in (11).

Algorithm 1: (NS-)DASF algorithm in fully-connected networks. Implementation available in [34]

```

begin
   $i \leftarrow 0, q \leftarrow 1$ , Randomly initialize  $X^0$ 
  loop
    for  $k \in \mathcal{K} \setminus \{q\}$  do
      At node  $k$ 
        Collect a new batch of  $N$  samples of
         $\mathbf{y}_k(t)$  and send the compressed
        samples  $\tilde{\mathbf{y}}_k^i(t) = X_k^{iT} \mathbf{y}_k(t)$  along with
         $\tilde{A}_k^i = X_k^{iT} A_k$ ,  $\tilde{B}_k^i = X_k^{iT} \tilde{B}_k$  and
         $\tilde{D}_k^i = X_k^{iT} D_k$  to node  $q$ .
      At node  $q$ 
        Obtain  $\tilde{X}^*$  by solving and selecting any
        solution of  $\mathbb{P}(\tilde{\mathbf{y}}^i, \tilde{B}^i, \tilde{D}^i, \tilde{A}^i)$ .
        Extract  $X_q^*$  and the  $\tilde{X}_k^*$ 's from  $\tilde{X}^*$ 
        according to (13).
         $X_q^{i+1} \leftarrow X_q^*$ 
        for  $k \in \mathcal{K} \setminus \{q\}$  do
          Send  $\tilde{X}_k^*$  to node  $k$ .
          At node  $k$ 
             $X_k^{i+1} \leftarrow X_k^i \tilde{X}_k^*$ 
         $i \leftarrow i + 1, q \leftarrow (q + 1) \bmod K$ 

```

computations or too much memory. In addition, an exact global solver for the problem of interest might not even exist. For example, many solvers for non-convex problems often produce a stationary point or local optimum rather than a global optimum. In this section, we show how the local solution step can be significantly relaxed to allow for common iterative solvers, not necessarily converging to optimal values, while still guaranteeing convergence to a stationary point of the centralized problem (1) or (2).

In what follows, we express (1)-(2) in the equivalent form

$$\min_X L(X) \quad (17)$$

where

$$L(X) \triangleq \varphi(X^T \mathbf{y}(t), X^T B) + \gamma(X^T A) + \delta_{\mathcal{X}}(X) \quad (18)$$

with $\delta_{\mathcal{X}}(\cdot)$ the indicator function of some set \mathcal{X} described by the equality and inequality constraints, i.e., $\delta_{\mathcal{X}}(X) = 0$ if $X \in \mathcal{X}$ and $\delta_{\mathcal{X}}(X) = \infty$ otherwise⁵. In the case of the generic problem (1) we have

$$\mathcal{X} \triangleq \{X \mid \forall j \in \mathcal{J}_I, \eta_j(X^T \mathbf{y}(t), X^T D_j) \leq 0, \forall j \in \mathcal{J}_E, \eta_j(X^T \mathbf{y}(t), X^T D_j) = 0\}, \quad (19)$$

and in the non-smooth case (2),

$$\mathcal{X} \triangleq \{X \mid \forall k, \forall j \in \mathcal{J}_I^k, \eta_j(X_k^T \mathbf{y}_k(t), X_k^T D_{j,k}) \leq 0, \forall j \in \mathcal{J}_E^k, \eta_j(X_k^T \mathbf{y}_k(t), X_k^T D_{j,k}) = 0\}. \quad (20)$$

⁵ L is thus not a real-valued function, but takes its values in the extended real number line $\mathbb{R} \triangleq \mathbb{R} \cup \{-\infty, \infty\}$.

From the update rules (14) and (15), one could show that solving $\mathbb{P}(\tilde{\mathbf{y}}^i, \tilde{B}^i, \tilde{D}^i, \tilde{A}^i)$ in the space where \tilde{X} lives is equivalent to solving $\mathbb{P}(\mathbf{y}, B, D, A)$ in the original space of X with the additional constraints that

$$\forall k \neq q, \exists \tilde{X}_k \in \mathbb{R}^{Q \times Q} : X_k = X_k^i \tilde{X}_k, \quad (21)$$

or equivalently

$$\forall k \neq q, X_k \in \text{range } X_k^i, \quad (22)$$

with range denoting the range or column space operator. To each local \tilde{X} is associated a unique X in the original problem space, a fact that we can express as the linear relationship

$$X = C_q^i \tilde{X} \quad (23)$$

where $C_q^i \in \mathbb{R}^{M \times (M_q + (K-1)Q)}$ is a particular block-matrix performing the required permutations and multiplications by the blocks X_k^i of the blocks \tilde{X}_k of \tilde{X} to obtain the transformation uniquely described by (21). More specifically, we define

$$\Theta_{-q}(X) \triangleq \text{BlockDiag}(X_1, \dots, X_{q-1}) \text{ and,} \quad (24)$$

$$\Theta_{+q}(X) \triangleq \text{BlockDiag}(X_{q+1}, \dots, X_K)$$

from which we define the linear matrix-to-matrix map

$$C_q(X) = \begin{bmatrix} O & \Theta_{-q}(X) & O \\ I_{M_q} & O & O \\ O & O & \Theta_{+q}(X) \end{bmatrix} \quad (25)$$

where the O 's denote all-zero matrices of appropriate dimensions. This allows us to formally express

$$C_q^i \triangleq C_q(X^i), \quad (26)$$

leading in particular to

$$X^i = C_q^i [X_q^{iT}, I_Q, \dots, I_Q]^T. \quad (27)$$

It is important to note that there is a linear relationship between the *local* variable \tilde{X} and the corresponding point X in the *global* space, and that this relationship depends on the current iteration i (via the previous estimate of the solution X^i) and the current updating node q . This relationship *itself* linearly depends on X^i via (26). The linearity, and hence continuity of the map C_q is one of the key properties used to show convergence, and replacing C_q by any other continuous map would lead to the same convergence result (although optimality would not anymore be guaranteed). Although we haven't described the algorithm for arbitrary network topologies, a similar linear map as in (25) can be defined for that case (see [28], [31]), which implies that the convergence proof in Section IV-C will generalize to these arbitrary topologies as well.

Using the notation (18) and (23), we can now compactly express the local problem at iteration i as

$$\min_{\tilde{X}} L(C_q^i \tilde{X}). \quad (28)$$

We note that –by construction– (28) has the same solutions as the local problem $\mathbb{P}(\tilde{\mathbf{y}}^i, \tilde{B}^i, \tilde{D}^i, \tilde{A}^i)$. While we need (28) in the proof, in practice it is better to directly solve $\mathbb{P}(\tilde{\mathbf{y}}^i, \tilde{B}^i, \tilde{D}^i, \tilde{A}^i)$

instead of (28) as the former directly inherits the same structure as the original centralized problem (1)-(2), and hence the original solver for (1)-(2) can be re-used.

A. Useful concepts from variational analysis

We introduce here a few definitions from variational analysis that are essential to our next developments. This section is supposedly kept short, and we refer the reader to the reference works of Clarke [35] and Rockafellar [36] for an in-depth treatment of the subject.

a) *Subgradients*: Following [36, Definition 8.3], we define the set of regular subgradients of L at \bar{X} (where $L(\bar{X})$ is finite) as

$$\hat{\partial}L(\bar{X}) \triangleq \left\{ V \in \mathbb{R}^{M \times Q} \mid \liminf_{\substack{X \rightarrow \bar{X} \\ X \neq \bar{X}}} \frac{L(X) - L(\bar{X}) - \langle V, X - \bar{X} \rangle_F}{\|X - \bar{X}\|} \geq 0 \right\}, \quad (29)$$

and the set of general subgradients of L at \bar{X} as⁶

$$\partial L(\bar{X}) \triangleq \{ V \in \mathbb{R}^{M \times Q} \mid \exists X^j \rightarrow \bar{X}, L(X^j) \rightarrow L(\bar{X}), V^j \in \hat{\partial}L(X^j) \rightarrow V \}, \quad (30)$$

which we refer to simply as “subgradients” from hereon. This definition of a subgradient does not require any structural property such as convexity, smoothness, or even continuity from L , and extends the usual definition of subgradients from convex analysis. Indeed, (30) generalizes the familiar notion of subgradient of a convex function. Note again that the developments in [31] assumed convexity for the non-smooth part of the objective in (2), which is not required here.

b) *Local Optimality*: If $L : \mathbb{R}^{M \times Q} \mapsto \mathbb{R}$ is lower-semicontinuous⁷ and X^* is a local minimum of L , then [36, Theorem 8.15]

$$0 \in \partial L(X^*). \quad (31)$$

Any point X^* such that $L(X^*)$ is finite and satisfying (31) is called a *stationary point* of L . This definition reduces to the usual “KKT” optimality conditions [37], [38] in the smooth case. Indeed, if

$$L(X) \triangleq f(X) + \delta_{\mathcal{X}}(X) \quad (32)$$

where f is a smooth function, and \mathcal{X} is the closed set $\{X \mid g(X) \leq 0, h(X) = 0\}$ where h and g are smooth functions, then we have at any point X where $L(X)$ is finite (and thus $X \in \mathcal{X}$) [39, Proposition 4.58]

$$\partial L(X) = \nabla f(X) + \partial \delta_{\mathcal{X}}(X). \quad (33)$$

If $\nabla g(X)$ and $\nabla h(X)$ are linearly independent (this is stricter than required, but is only meant for illustrative purposes), then [39, Example 4.49]

$$\partial \delta_{\mathcal{X}}(X) = \{ \lambda_g \nabla g(X) + \lambda_h \nabla h(X) \mid \lambda_h \in \mathbb{R}, \lambda_g = 0 \text{ if } g(X) < 0, \text{ else } \lambda_g \geq 0 \}. \quad (34)$$

⁶(30) can be read as “the set of elements V such that there exist sequences $(X^j)_j$ converging to \bar{X} , $(L(X^j))_j$ converging to $L(\bar{X})$ and $(V^j)_j$ converging to V , such that V^j is a regular subgradient of L at X^j ”.

⁷A function L is lower-semicontinuous if for any α , its level sets $\{X \mid L(X) \leq \alpha\}$ are closed.

Defining the *Lagrangian*

$$\mathfrak{L}(X, \lambda_g, \lambda_h) \triangleq f(X) + \lambda_g g(X) + \lambda_h h(X), \quad (35)$$

we have for a feasible point X the equivalence

$$0 \in \partial L(X) \Leftrightarrow \begin{aligned} &\exists \lambda_g = 0 \text{ if } g(X) < 0, \text{ else } \lambda_g \geq 0, \lambda_h : \\ &\quad \nabla_X \mathfrak{L}(X, \lambda_g, \lambda_h) = 0. \end{aligned} \quad (36)$$

This equivalence allows us to offer a unifying treatment of the smooth and non-smooth cases via the compact problem formulation (17) (or (28) for the local problem at the updating node).

B. DASF algorithm with inexact solver

In what follows, we define the local objective at iteration i

$$\tilde{L}_q^i : \tilde{X} \mapsto L(C_q^i \tilde{X}), \quad (37)$$

which we use to construct the following definition of an *inexact iterative solver*:

Definition 1 (Inexact iterative solver). *Given some current estimate X^i , an inexact iterative local solver running on node q produces a sequence of iterates $(\tilde{X}^{i,j})_{j \in \mathbb{N}}$ such that for each iteration i*

- 1)
$$\tilde{X}^{i,0} = [X_q^{iT}, I_Q, \dots, I_Q]^T, \quad (C1)$$

- 2) *there is some $R^i > 0$ such that for any j ,*

$$\begin{aligned} &\tilde{L}_q^i(\tilde{X}^{i,j}) - \tilde{L}_q^i(\tilde{X}^{i,j+1}) \\ &\geq R^i \left\| \tilde{X}^{i,j} - \tilde{X}^{i,j+1} \right\|_F^2, \end{aligned} \quad (C2)$$

- 3) *there is some $c^i \geq 0$ such that for any j there is some $W \in \partial \tilde{L}_q^i(\tilde{X}^{i,j+1})$ such that*

$$c^i \left\| \tilde{X}^{i,j} - \tilde{X}^{i,j+1} \right\|_F \geq \|W\|_F. \quad (C3)$$

- 4) *In addition, the sequences of constants R^i and c^i must be such that*

$$\liminf_{i \rightarrow \infty} R^i > 0 \quad \text{and} \quad \limsup_{i \rightarrow \infty} c^i < \infty. \quad (C4)$$

Flowing directly from (27), condition (C1) ensures consistency with the current global estimate of the solution X^i , and can be enforced as long as the solver accepts an arbitrarily chosen starting point. Condition (C2) is often referred to as a *sufficient decrease* condition [40], [41]. It ensures that an update of the optimization variable results in a corresponding decrease of the cost, and produces a feasible point (as, by the definition of the indicator function $\delta_{\mathcal{X}}$, the cost would otherwise have infinite value in the case of a non-feasible point). Condition (C3) ensures that progressively smaller updates of the optimization variable result in a higher degree of stationarity (the distance of the set $\partial \tilde{L}(\tilde{X}^{j+1})$ to 0 plays here the role of optimality criterion, and summarizes how “close” we are to a stationary point). This third condition is crucial, as

without it, the procedure could wander away from a stationary point during the first few steps, but still eventually converge to some stationary point, in which case a single step of the solver would not be guaranteed to improve on the current iterate. The last requirements on R^i and c^i in condition (C4) ensure that those constants do not become arbitrarily small or large as the algorithm progresses, as that would render the two previous conditions meaningless. Those constants are usually related to a parameter of the solver that can be controlled, a lower or upper bound can therefore be enforced explicitly.

Although these requirements might seem restrictive, they have been shown to hold for many commonly used descent methods [42], if some additional technical conditions are satisfied. Here are some common examples, some of which are further elaborated on in Appendix A:

Gradient Descent (see Appendix A, and [42], [43]) provided that the problem is unconstrained and the objective has Lipschitz gradients.

Projected/Proximal Gradient Descent (see [40], [42]) provided that the smooth part of the objective has Lipschitz gradients.

Newton's Method (see Appendix A) provided that the Hessian is positive definite and bounded, and that the objective has Lipschitz gradients.

Regularized Exact Solver (see Appendix A) By definition, an exact solver obtains a solution in a single iteration and the properties of Definition 1 are almost satisfied. If there are multiple solutions, some mechanism is required to ensure that the solver does not “jump” from one solution to the next, by for example explicitly penalizing the distance from the starting point.

Power Method (see Appendix A) An intuitive argument is that the power method can be expressed as an instance of projected gradient descent applied to a specific constrained quadratic problem.

Remark IV.1 (Comparison with the original DASF algorithm). An important, mostly theoretical, difference with the original DASF and NS-DASF algorithms, is the removal of the *well-posedness* assumption on the global problems. Indeed, it was assumed in those algorithms that the solution set of the optimization problem varied “smoothly” with regards to the problem's parameters $(\mathbf{y}(t), B, \mathcal{D}, A)$. This assumption can be hard to check, and even not hold in some cases. Property (C2), which is not satisfied by an exact solver, allows us to remove this assumption. See [28], [29], [31] for details on the assumption.

The inexact version of the DASF Algorithm simply consists in selecting $\tilde{X}^* \triangleq \tilde{X}^{i, n_i}$ as the solution of an inexact solver applied to the local problem at node q at iteration i of the DASF algorithm, which we denote

$$\min_{\tilde{X}} \tilde{L}_q^i(\tilde{X}), \quad (38)$$

where the number of iterations $n_i > 0$ of the inexact solver can be chosen arbitrarily.

C. Convergence and Optimality

We first describe a general and easily satisfied assumption on L .

Assumption 1. L is continuous over its domain⁸, and it has compact sublevel-sets.

A sufficient condition is to require that φ and the η_j are continuous, γ is continuous over its domain, and either γ and φ or at least one of the η_j have bounded sublevel sets.

We show the convergence of this scheme by showing that (i) $L(X^i)$ decreases monotonically, and (ii) every accumulation point of $(X^i)_{i \in \mathbb{N}}$ is a stationary point (as defined by (31)) of the local problem associated with each node q (38). We then leverage the main result of [29], [31] to show optimality with regards to the global optimization problem (17).

Lemma 1 (Monotonic decrease of the objective). *Let $(X^i)_{i \in \mathbb{N}}$ be a sequence generated by the DASF Algorithm with inexact local solver. Then $(L(X^i))_{i \in \mathbb{N}}$ is a monotonically decreasing convergent sequence.*

Proof. Since by the definition (37) $\tilde{L}_q^i(\tilde{X}) = L(C_q^i \tilde{X})$, from the second property of the local solver (C2), we have

$$L(C_q^i \tilde{X}^{i,j}) - L(C_q^i \tilde{X}^{i,j+1}) \geq R^i \left\| \tilde{X}^{i,j} - \tilde{X}^{i,j+1} \right\|_F^2 \geq 0. \quad (39)$$

Summing from $j = 0$ to $j = n_i$ and telescoping the left-hand side yields

$$L(X^i) - L(X^{i+1}) \geq 0, \quad (40)$$

where we used (C1) to substitute $X^i = C_q^i \tilde{X}^{i,0}$ and the definition of the inexact algorithm to substitute $X^{i+1} = C_q^i \tilde{X}^{i, n_i}$, hence proving the monotonic decrease.

As L is continuous over its domain with compact sub-level sets, it has a minimum value, and the sequence $(L(X^i))_{i \in \mathbb{N}}$ is therefore bounded-below [44] and must converge [45]. \square

Lemma 2 (Convergence of the Sequence of Residuals). *Let $(X^i)_{i \in \mathbb{N}}$ be a sequence generated by the DASF Algorithm with inexact local solver, and let $(\tilde{X}^{i,j})_{j \in \mathbb{N}}$ be the subsequence generated by the local solver at each iteration i . Then ,*

$$\lim_{i \rightarrow \infty} \left\| \tilde{X}^{i,j} - \tilde{X}^{i,j+1} \right\|_F^2 = 0 \quad \forall j \quad (41)$$

and

$$\lim_{i \rightarrow \infty} \left\| X^i - X^{i+1} \right\|_F = 0. \quad (42)$$

Proof. Using the fact that

$$X^i = C_q^i \tilde{X}^{i,0} \text{ and } X^{i+1} = C_q^i \tilde{X}^{i, n_i}, \quad (43)$$

we have

$$L(X^i) - L(X^{i+1}) = \sum_{j=0}^{j=n_i} L(C_q^i \tilde{X}^{i,j}) - L(C_q^i \tilde{X}^{i,j+1}), \quad (44)$$

⁸The domain of a function L is the set of points for which $L(X) < \infty$.

From the second property of the local solver (C2), for any i and $j \leq n_i$

$$L(C_q^i \tilde{X}^{i,j}) - L(C_q^i \tilde{X}^{i,j+1}) \geq 0, \quad (45)$$

it must be that, for any i and $j \leq n_i$

$$L(X^i) - L(X^{i+1}) \geq L(C_q^i \tilde{X}^{i,j}) - L(C_q^i \tilde{X}^{i,j+1}), \quad (46)$$

and hence from (C2) again

$$L(X^i) - L(X^{i+1}) \geq R^i \left\| \tilde{X}^{i,j} - \tilde{X}^{i,j+1} \right\|_F^2. \quad (47)$$

As a consequence of Lemma 1, the convergence of $(L(X^i))_{i \in \mathbb{N}}$ in combination with (47) implies that, for any $0 \leq j < n_i$,

$$\lim_{i \rightarrow \infty} L(X^i) - L(X^{i+1}) = \lim_{i \rightarrow \infty} R^i \left\| \tilde{X}^{i,j} - \tilde{X}^{i,j+1} \right\|_F^2 = 0. \quad (48)$$

From property (C4) of the local solver, the sequence R^i is eventually lower-bounded, and therefore

$$\lim_{i \rightarrow \infty} \left\| \tilde{X}^{i,j} - \tilde{X}^{i,j+1} \right\|_F^2 = 0 \quad \forall j, \quad (49)$$

proving the first statement.

From (43) and the triangle inequality, we have

$$\begin{aligned} \|X^i - X^{i+1}\|_F &\leq \sum_{j=0}^{n_i-1} \left\| C_q^i \tilde{X}^{i,j} - C_q^i \tilde{X}^{i,j+1} \right\|_F \leq \\ &\|C_q^i\|_F \sum_{j=0}^{n_i-1} \left\| \tilde{X}^{i,j} - \tilde{X}^{i,j+1} \right\|_F, \end{aligned} \quad (50)$$

where the second inequality relies on the sub-multiplicative property of the Frobenius norm. As $C_q^i = C_q(X^i)$, and as the sub-level sets of $L(X^0)$ are compact, $\|X^i\|_F$ is bounded and by the continuity of the map $C_q(X)$, so is $\|C_q^i\|$. (49) and (50) therefore imply that

$$\lim_{i \rightarrow \infty} \|X^i - X^{i+1}\|_F = 0, \quad (51)$$

completing the proof. \square

Lemma 3 (Node-specific Optimality of Accumulation Points). *Let $(X^i)_{i \in \mathbb{N}}$ be a sequence generated by the DASF Algorithm with inexact local solver. Then, every accumulation point \bar{X} of $(X^i)_{i \in \mathbb{N}}$ is a stationary point of the local problem (38) for any updating node q .*

Proof. By the third and fourth property of the local solver (C3)-(C4), we can find sequences $(W^i)_{i \in \mathbb{N}}$ and $(c^i)_{i \in \mathbb{N}}$ (with $c^i > 0$) such that

$$W^i \in \partial \tilde{L}_q^i(\tilde{X}^{i,n_i}), \quad (52)$$

and

$$c \left\| \tilde{X}^{i,n_i-1} - \tilde{X}^{i,n_i} \right\|_F \geq \|W^i\|_F, \quad (53)$$

where $c = \limsup_i c^i$. We can express equation (52) in the global domain of X by applying the generalized chain rule^{9,10} [36] to $L(C_q^i \tilde{X}^{i,n_i})$:

$$\begin{aligned} \partial \tilde{L}_q^i(\tilde{X}^{i,n_i}) &= \partial_{\tilde{X}} \left(L(C_q^i \tilde{X}^{i,n_i}) \right) \\ &= C_q^{iT} \partial L(C_q^i \tilde{X}^{i,n_i}) = C_q^{iT} \partial L(X^{i+1}) \end{aligned} \quad (54)$$

where we used the fact that $C_q^i \tilde{X}^{i,n_i} = X^{i+1}$ by definition of the algorithm.

If \bar{X} is an accumulation point, then there is (by definition) some index set \mathcal{I} such that¹¹ $\lim_{i \in \mathcal{I} \rightarrow \infty} X^i = \bar{X}$. Furthermore, from Lemma 2 and in particular (42), we also have

$$\lim_{i \in \mathcal{I} \rightarrow \infty} X^{i+1} = \lim_{i \in \mathcal{I} \rightarrow \infty} X^i = \bar{X}. \quad (55)$$

As the number of nodes is finite, we can furthermore select \mathcal{I} such that the sequence of updating nodes $(q^i)_{i \in \mathcal{I}}$ is a constant sequence, i.e. we only consider a single updating node q . From the continuity of the map $C_q(X)$ it must also be that there is some $\bar{C}_q \triangleq C_q(\bar{X})$ such that

$$\lim_{i \in \mathcal{I} \rightarrow \infty} C_q^i = \lim_{i \in \mathcal{I} \rightarrow \infty} C_q(X^i) = C_q(\bar{X}) = \bar{C}_q. \quad (56)$$

From (52), (54) and the outer-semicontinuity of the set of subgradients [36], we have

$$W^i \in C_q^{iT} \partial L(X^{i+1}) \forall i \Rightarrow \lim_{i \in \mathcal{I} \rightarrow \infty} W^i \in \bar{C}_q^T \partial L(\bar{X}). \quad (57)$$

From (41) and (53), it must be that

$$\lim_{i \rightarrow \infty} \|W^i\|_F = 0, \quad (58)$$

and thus

$$\lim_{i \in \mathcal{I} \rightarrow \infty} W^i = 0. \quad (59)$$

Combining (59) with (57) yields

$$0 \in C_q(\bar{X})^T \partial L(\bar{X}) \quad (60)$$

which is the condition for \bar{X} to be a stationary point of the local problem (38), based on the argument in (54). Note that the local problem (38) is equal to (17) equipped with the additional constraint (21), i.e. \bar{X} is a stationary point of

$$\min_X L(X) \quad \text{s.t.} \quad X = C_q(\bar{X}) \tilde{X}, \quad \tilde{X} \in \mathbb{R}^{((K-1)Q+M_q) \times Q}.$$

We have shown that an accumulation point is a stationary point for at least a single node q , it remains to show that it is the case for any q . From Lemma 2 and in particular (42), if \bar{X} is an accumulation point of $(X^{i+1})_{i \in \mathcal{I}}$, then it is also an accumulation point of $(X^{i+n})_{i \in \mathcal{I}}$ for any $n \geq 0$. Indeed,

$$X^{i+n} = X^i - (X^i - X^{i+1}) - \dots - (X^{i+n-1} - X^{i+n}), \quad (61)$$

⁹We denote $\partial_{\tilde{X}}(\cdot)$ the set of subgradients of the expression in argument, interpreted as a function of the local variable \tilde{X} .

¹⁰We define the product between a set and a matrix as the set whose elements are the elements of the original set multiplied by that matrix.

¹¹ $\lim_{i \in \mathcal{I} \rightarrow \infty} a_i$ is a short-hand notation for $\lim_{j \rightarrow \infty} a_{i_j}$ where i_j is the j -th element in the ordered index set \mathcal{I} .

and therefore

$$\begin{aligned} \lim_{i \rightarrow \infty} X^{i+n} - X^i = \\ \lim_{i \rightarrow \infty} (X^i - X^{i+1}) - \dots - (X^{i+n-1} - X^{i+n}) = 0, \end{aligned} \quad (62)$$

and finally

$$\lim_{i \in \mathcal{I} \rightarrow \infty} X^i = \lim_{i \in \mathcal{I} \rightarrow \infty} X^{i+n} = \bar{X}. \quad (63)$$

As we have shown that \bar{X} is a stationary point for at least one node, it is a stationary point for all nodes as the sequence $(X^{i+n})_{i \in \mathcal{I}}$ will be associated with node¹² $(q+n) \bmod K$. \square

Lemma 3 asserts that any accumulation point of DASF is a stationary point of all local problems at all nodes. We will use this result to show that these accumulation points are also stationary points of the centralized problem (1)-(2). For this, we need a technical condition that is usually satisfied with high probability if the number of constraints is not too high (see below). The condition can actually be viewed as a compressed version of the standard linear independence constraint qualifications (LICQ) in numerical optimization [43], except that they apply to the ‘‘compressed’’ gradients rather than the actual gradients of the constraint functions.

Proposition 1 (Compressed LICQ). *Let $\vartheta_j : X \mapsto \eta_j(X^T \mathbf{g}, X^T D_j)$ in the smooth case, and $\vartheta_j : X \mapsto \eta_j(X_k^T \mathbf{g}_k, X_k^T D_{j,k})$ in the non-smooth case, and define the active constraint set $\mathcal{A}(X) \triangleq \{j \in \mathcal{I}_I \mid \vartheta_j(X) = 0\}$. If the elements of the set*

$$\{\text{BlockDiag}(\bar{X}_1, \dots, \bar{X}_K)^T \nabla \vartheta_j(\bar{X}) \mid j \in \mathcal{A}(\bar{X}) \cup \mathcal{J}_E\} \quad (64)$$

are linearly independent matrices and \bar{X} is a stationary point of the local problems for any node q , i.e.

$$\forall q \in \mathcal{K}, 0 \in \bar{C}_q^T \partial L(\bar{X}), \quad (65)$$

then it is also a stationary point of the global problem (1)-(2), i.e.

$$0 \in \partial L(\bar{X}). \quad (66)$$

A proof of the above statement is available in [29] for the smooth case and in [31] for the non-smooth case. Note that this qualification is automatically violated if the number of active constraints exceeds KQ^2 (which is the dimension of the vector space in which the compressed gradients live). See [28], [29] for details and an extension of this condition to the case of arbitrary network topologies.

¹²assuming a sequential selection rule, but the conclusion stays valid as long as every node is selected infinitely many times.

Theorem 1 (Convergence and Optimality). *Let $(X^i)_{i \in \mathbb{N}}$ be a sequence generated by the DASF algorithm with inexact solver. Then if Assumption 1 is satisfied and the qualification (64) is also satisfied at the accumulation points of $(X^i)_{i \in \mathbb{N}}$, $(X^i)_{i \in \mathbb{N}}$ converges to the set of stationary points of the global problem (1)-(2), that is*

$$\lim_{i \rightarrow \infty} \min_{X \in \Omega} \|X - X^i\|_F = 0, \quad (67)$$

where Ω denotes the set of stationary points of problem (1)-(2).

Furthermore, if the number of stationary points of (1)-(2) is finite, then $(X^i)_{i \in \mathbb{N}}$ converges to a single point.

Proof. Under Assumption 1, Lemma 3 asserts that any accumulation point is a stationary point of the local problem at every node. Proposition 1 ensures that under the qualification (64), the local stationarity at every node implies stationarity for the global problem, proving the first part of the theorem (i.e., convergence to the set of stationary points). The second part (convergence to a single point) is a direct consequence of Lemma 2 and is a standard analysis result. See the proof of [29, Theorem 4] for details. \square

As a consequence of the above theorem, interleaving the steps of DASF with a single descent step is sufficient to ensure convergence to a stationary point.

D. A Bound on the Convergence Rate

Obtaining a convergence rate for the objective value would unfortunately require much more assumptions on the functions involved in (1)-(2), than what we assumed so far. Still, we can obtain a lower bound on the convergence rate of a local optimality measure at each node, that is

$$w^i \triangleq \text{dist}(0, \partial \tilde{L}_{q^i}^i(\tilde{X}^{i,n_i})) \triangleq \min_{W \in \partial \tilde{L}_{q^i}^i(\tilde{X}^{i,n_i})} \|W\|. \quad (68)$$

As our actual objective is to find stationary points rather than an actual minimum, w^i is an appropriate measure of optimality. We could also have considered the step-length $\|X^i - X^{i+1}\|$ as a measure of optimality, which would yield a rate similar to the one presented hereafter.

Proposition 2 (Convergence Rate Bound). *Let $(X^i)_{i \in \mathbb{N}}$ be a sequence generated by the DASF algorithm with inexact solver, and let w^i be as defined in (68). Then there is some positive constant a such that*

$$\min_{j \leq i} w^j \leq a \frac{\sqrt{L(X^0) - L^*}}{\sqrt{i+1}}, \quad (69)$$

where L^ is the minimum value of L .*

Proof. This follows from a well-known proof argument [40], [41]. From (47) and condition (C3), we have

$$L(X^i) - L(X^{i+1}) \geq \frac{R^i}{c^i} (w^i)^2. \quad (70)$$

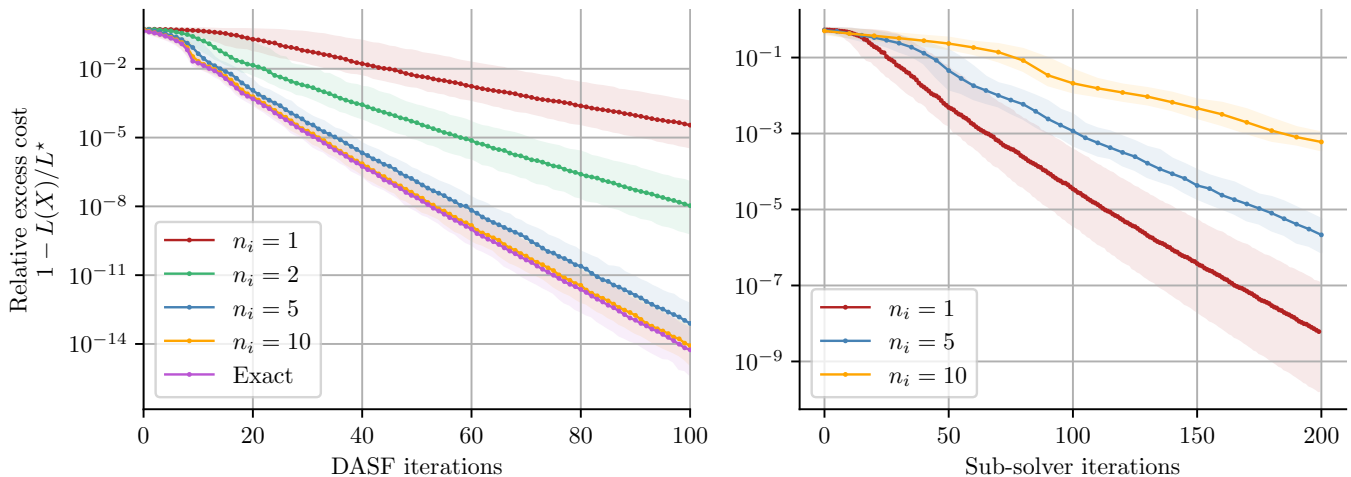


Fig. 2: Convergence of the DASF algorithm with inexact solver for the Max-SNR problem. Each curve was generated by 1000 Monte-Carlo runs. Solid curves depict the median trajectories, while the shaded areas of corresponding color depict the the 5%-95% percentiles regions.

Summing the inequality from 0 to i and telescoping yields

$$L(X^0) - L(X^{i+1}) \geq \sum_{j=0}^i \frac{R^j}{c^j} (w^j)^2. \quad (71)$$

From condition (C4), there is some $r > 0$ such that

$$\liminf_{i \rightarrow \infty} \frac{R^i}{c^i} = r, \quad (72)$$

and hence, using the fact that $L^* \leq L(X^{i+1})$,

$$L(X^0) - L^* \geq r \sum_{j=0}^i (w^j)^2 \geq r(i+1) \min_j (w_j)^2, \quad (73)$$

which can be rearranged to obtain (69) with $a = r^{-1/2}$. \square

For a given level of optimality $\min w^i$, the bound gives a linear relationship between the required number of iterations and the current error $L(X^i) - L^*$. In other words, if the algorithm is set to stop after a sufficiently small norm of the subgradient is reached, we can guarantee that the algorithm will run for a number of iterations (hence time) that is at most proportional to the current error $L(X^i) - L^*$.

V. NUMERICAL EXPERIMENTS

In order to qualitatively demonstrate the convergence of the DASF algorithm with inexact solver, we consider an SNR maximization (Max-SNR) problem [5] defined as

$$\max_x \mathbb{E} \{ \|x^T \mathbf{y}(t)\|^2 \} \quad \text{s.t.} \quad \mathbb{E} \{ \|x^T \mathbf{n}(t)\|_F^2 \} = 1 \quad (74)$$

where x denotes the $M \times 1$ spatial filter ($Q = 1$), $\mathbf{y}(t)$ is modelled as $\mathbf{y}(t) = \mathbf{a}d(t) + \mathbf{n}(t)$, $\mathbf{n}(t)$ is an M -dimensional i.i.d white Gaussian noise signal, $d(t)$ is a single-channel (unknown) target signal, and where \mathbf{a} is some unknown mixing vector with entries sampled from $\mathcal{N}(0, 1)$. The purpose of the multi-input single-output (MISO) filter x is to maximize

the SNR of d in the output filtered signal $z(t) = x^T \mathbf{y}(t)$. It is assumed that the nodes of the WSN are able to collect samples of both $\mathbf{y}(t)$ and $\mathbf{n}(t)$. For our simulations we set $M = 100$, $M_k = 10$, $K = 10$, i.e., there are 10 nodes, with each node collecting 10 different channels of the 100-channel signals $\mathbf{y}(t)$ and $\mathbf{n}(t)$. The noise variance is 10, while the signal variance is 1. Our figure of merit is the relative excess cost, defined as

$$1 - \frac{L(X)}{L^*}, \quad (75)$$

where L^* denotes the optimal objective value of (74). Note that the solution of (74) can be solved by computing a generalized eigenvalue decomposition (GEVD) on the covariance matrices $\mathbb{E} \{ \mathbf{y}(t) \mathbf{y}(t)^T \}$ and $\mathbb{E} \{ \mathbf{n}(t) \mathbf{n}(t)^T \}$ [46].

Figure 2 depicts the convergence of the inexact version of DASF applied to problem (74) with different numbers of sub-solver iterations (n^i). The exact (GEVD) solver uses Cuppen's divide and conquer algorithm [47], while the inexact solver relies on the generalized power method [48]. On the left part of the figure, we see that the algorithm converges to the optimal objective value with a single sub-solver iteration per-node, and that 10 sub-solver iterations are sufficient to reach a convergence rate at par with the exact solver. Although, based on this first plot, one might conclude that more sub-solver iterations is better, the right part of Figure 2 tells a very different story. The same figure of merit is depicted, but in terms of the total number of sub-solver iterations, rather than "global" DASF iterations. This second plot therefore depicts the algorithm's convergence as a function of "real" time and offers a fairer comparison. Indeed, the sub-solver with the lowest number of iterations ($n_i = 1$) outperforms all the other ones. It seems that the smaller progress achieved at each step is largely compensated by the fact that the rate at which the updating node role is passed-on is much larger, resulting in the network being able to update the filter in more "directions" (i.e. with different subspace constraints (21)) in a given time-interval, albeit with smaller progress at each updating node.

These more frequent updates of the updating node do come with a cost, as a new batch of N compressed samples data should normally be retransmitted every time the updating node changes. This can be largely compensated by applying the data-reuse technique described in [49], making the DASF algorithm with inexact local solver an excellent option for adaptively tracking a spatial filter.

VI. CONCLUSION

We have extended the convergence results of the DASF algorithm to a setting where the exact solver is replaced by an inexact, iterative one. We have shown, both by the means of proofs and simulations, that the convergence properties of DASF were preserved under more relaxed technical conditions. Furthermore, the simulations have shown that DASF used with an iterative solver could display better convergence properties, at the cost of a slightly increased bandwidth, assuming a data-reuse scheme as in [49].

APPENDIX

A. Common Inexact Solvers

We describe in this section a couple of methods satisfying the requirements set by Definition 1. In what follows, we consider methods to minimize a real-valued function $h : \mathbb{R}^{M \times Q} \mapsto \mathbb{R}$, whose precise structure (e.g. smoothness, convexity, etc.) will be redefined for every method. We first state a useful definition:

Definition 2. A smooth function $h : \mathbb{R}^{M \times Q} \mapsto \mathbb{R}$ is said to have R -Lipschitz gradients if there is some $R > 0$ such that for any $X, Y \in \mathbb{R}^{M \times Q}$

$$\|\nabla h(X) - \nabla h(Y)\|_F \leq R \|X - Y\|_F. \quad (76)$$

In addition, [40, Descent Lemma], (76) also implies that

$$h(Y) \leq h(X) + \langle \nabla h(X), Y - X \rangle_F + \frac{R}{2} \|X - Y\|_F^2, \quad (77)$$

where $\langle A, B \rangle_F \triangleq \text{Tr}(A^T B)$.

1) *Line-Search Methods: Gradient Descent and Newton's Method:* Line search methods are defined by the update rule

$$X^{i+1} = X^i + \mu^i P^i, \quad (78)$$

where P^i is the *search direction* and μ^i is an appropriately chosen step-size.

a) *Gradient Descent:* In the case of Gradient Descent, we have $P^i = -\nabla h(X^i)$ and hence

$$X^{i+1} - X^i = -\mu^i \nabla h(X^i). \quad (79)$$

As a consequence, we have

$$\frac{1}{\inf_i \mu^i} \|X^{i+1} - X^i\| \geq \|\nabla h(X^i)\|. \quad (80)$$

As h is smooth, $\partial h(X^i) = \{\nabla h(X^i)\}$ [36], and condition (C3) of Definition 1 is therefore satisfied.

Applying the Descent Lemma (77) to X^i, X^{i+1} we have

$$\begin{aligned} & h(X^i) - h(X^{i+1}) \\ & \geq -\langle \nabla h(X^i), X^{i+1} - X^i \rangle_F - \frac{R}{2} \|X^{i+1} - X^i\|_F^2. \end{aligned} \quad (81)$$

From (79), (81) becomes

$$h(X^i) - h(X^{i+1}) \geq \left(\frac{1}{\mu^i} - \frac{R}{2} \right) \|X^{i+1} - X^i\|_F^2. \quad (82)$$

and condition (C2) is thus also satisfied if $\sup_i \mu^i < \frac{2}{R}$.

b) *Newton's Method:* Newton's step consists in setting $P^i = -(\nabla^2 h(X^i))^{-1} \nabla h(X^i)$. Hence, from the update rule (78), we have

$$-\frac{1}{\mu^i} \nabla^2 h(X^i) (X^{i+1} - X^i) = \nabla h(X^i). \quad (83)$$

Taking the norm on both sides and denoting as λ_{\max} an upper bound on the largest eigenvalue of the Hessian (across all iterations),

$$\frac{\lambda_{\max}}{\inf_i \mu^i} \|X^{i+1} - X^i\| \geq \|\nabla h(X^i)\|, \quad (84)$$

and condition (C3) is fulfilled by the same reasoning as for gradient descent.

Taking the inner product of both sides of (83) with $(X^{i+1} - X^i)$ and denoting the Cholesky factorization of the Hessian as $\nabla^2 f(X^i) \triangleq U U^T$, we have

$$-\langle \nabla h(X^i), X^{i+1} - X^i \rangle_F = \frac{1}{\mu^i} \|U^T (X^{i+1} - X^i)\|^2. \quad (85)$$

Hence

$$-\langle \nabla h(X^i), X^{i+1} - X^i \rangle_F \geq \frac{\lambda_{\min}}{\sup_i \mu^i} \|X^{i+1} - X^i\|^2 \quad (86)$$

where λ_{\min} denotes a lower bound on the smallest eigenvalue of the Hessian (across all iterations). Combining the above inequality with (81) yields

$$h(X^i) - h(X^{i+1}) \geq \left(\frac{\lambda_{\min}}{\sup_i \mu^i} - \frac{R}{2} \right) \|X^{i+1} - X^i\|_F^2, \quad (87)$$

and condition (C2) is satisfied as long as the step-size satisfies $\sup_i \mu^i < \frac{2\lambda_{\min}}{R}$ (note that this also implies that the Hessian must be positive definite).

In practice, the bounds on the step-size cannot be known apriori, but the same results can be obtained if a backtracking line-search [40], [43], [50] is used to obtain the step-size (we omit the analysis for brevity).

2) *Regularized Exact Solver:* Given some previous estimate of the solution X^- , a regularized exact solver would update the solution as (we drop any iteration index, as the solver only performs a single iteration)

$$X^+ = \underset{X}{\operatorname{argmin}} h(X) + \mu \|X - X^-\|_F^2 \quad (88)$$

instead of $\min_X h(X)$, where $\mu > 0$ is a freely chosen parameter. From the optimality of X^+ in (88), we have

$$h(X^-) - h(X^+) \geq \mu \|X^+ - X^-\|_F^2, \quad (89)$$

and condition (C2) is satisfied. Furthermore, the optimality of X^+ also implies that (using the sum-rule for subgradients [39])

$$0 \in \partial h(X^+) + 2\mu(X^+ - X^-), \quad (90)$$

And thus there is $W \triangleq 2\mu(X^- - X^+) \in \partial h(X^+)$, and condition (C3) is trivially satisfied.

3) *Power Method*: We will rely on the proof available in [42] that the projected gradient algorithm satisfies conditions (C2) and (C3) of Definition 1, even when the constraint set is non-convex. The projected gradient algorithm is defined by

$$X^{i+1} = P_{\mathcal{C}}(X^i - \mu \nabla h(X^i)) \quad (91)$$

where $P_{\mathcal{C}}$ is the projection on some constraint set \mathcal{C} and μ some step-size. We will show that the power method is a particular case of (91). The power method finds the eigenvector associated with the largest eigenvalue of some matrix A . It therefore finds a solution of

$$\min_x -x^T A x \quad \text{s.t.} \quad x^T x = 1. \quad (92)$$

The power method's update rule is

$$x^{i+1} = \frac{Ax^i}{\|Ax^i\|}. \quad (93)$$

Let us now assume that we solve the following (equivalent) problem using a projected gradient algorithm:

$$\min_x -x^T \frac{(A - I)}{2\mu} x \quad \text{s.t.} \quad x^T x = 1. \quad (94)$$

Denoting the objective h , its gradient is

$$\nabla h(x) = \frac{(I - A)}{\mu} x, \quad (95)$$

The projected gradient method with step-size μ for this problem is thus

$$x^{i+1} = P_{\mathcal{C}}\left(x^i - \mu \frac{x^i - Ax^i}{\mu}\right) = P_{\mathcal{C}}(Ax^i) \quad (96)$$

with \mathcal{C} denoting in this case the unit ball, which is equivalent to the update rule of the power method (93). Note that we skipped the discussion on the condition that the step-size should be smaller than the inverse of the Lipschitz constant of $\nabla h(x)$ [42], and simply mention that this condition can always be enforced by applying the proper scaling on A , as this changes neither the solution of (92) nor the update rule (93).

REFERENCES

- [1] H. Hotelling, "Analysis of a complex of statistical variables into principal components." *Journal of educational psychology*, vol. 24, no. 6, p. 417, 1933.
- [2] N. Wiener, *Extrapolation, interpolation, and smoothing of stationary time series: with engineering applications*. The MIT press, 1949.
- [3] M. S. Bartlett, "The statistical significance of canonical correlations," *Biometrika*, vol. 32, no. 1, pp. 29–37, 1941.
- [4] S. Markovich-Golan, A. Bertrand, M. Moonen, and S. Gannot, "Optimal distributed minimum-variance beamforming approaches for speech enhancement in wireless acoustic sensor networks," *Signal Processing*, vol. 107, pp. 4–20, 2015.
- [5] B. D. Van Veen and K. M. Buckley, "Beamforming: A versatile approach to spatial filtering," *IEEE assp magazine*, vol. 5, no. 2, pp. 4–24, 1988.
- [6] P. Corke, T. Wark, R. Jurdak, W. Hu, P. Valencia, and D. Moore, "Environmental wireless sensor networks," *Proceedings of the IEEE*, vol. 98, no. 11, pp. 1903–1917, 2010.
- [7] A. Bertrand, J. Callebaut, and M. Moonen, "Adaptive distributed noise reduction for speech enhancement in wireless acoustic sensor networks," in *Proc. of the International Workshop on Acoustic Echo and Noise Control (IWAENC)*, 2010.
- [8] O. Elharrouss, N. Almaadeed, and S. Al-Maadeed, "A review of video surveillance systems," *Journal of Visual Communication and Image Representation*, vol. 77, p. 103116, 2021.
- [9] A. M. Narayanan, P. Patrinos, and A. Bertrand, "Optimal versus approximate channel selection methods for EEG decoding with application to topology-constrained neuro-sensor networks," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 29, pp. 92–102, 2020.
- [10] A. M. Narayanan, R. Zink, and A. Bertrand, "EEG miniaturization limits for stimulus decoding with EEG sensor networks," *Journal of Neural Engineering*, vol. 18, no. 5, p. 056042, 2021.
- [11] M. Inggs, H. Griffiths, F. Fioranelli, M. Ritchie, and K. Woodbridge, "Multistatic radar: System requirements and experimental validation," in *2014 International Radar Conference*. IEEE, 2014, pp. 1–6.
- [12] M. Inggs and A. van der Byl, "NeXtRAD and RHINO radar: Harnessing the herd for networked radar," in *2014 International Radar Conference*. IEEE, 2014, pp. 1–5.
- [13] S. Haykin and K. R. Liu, *Handbook on array processing and sensor networks*. John Wiley & Sons, 2010, vol. 63.
- [14] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein *et al.*, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends® in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [15] C. G. Lopes and A. H. Sayed, "Diffusion least-mean squares over adaptive networks: Formulation and performance analysis," *IEEE Transactions on Signal Processing*, vol. 56, no. 7, pp. 3122–3136, 2008.
- [16] A. H. Sayed *et al.*, "Adaptation, learning, and optimization over networks," *Foundations and Trends in Machine Learning*, vol. 7, no. 4–5, pp. 311–801, 2014.
- [17] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE signal processing magazine*, vol. 37, no. 3, pp. 50–60, 2020.
- [18] C. I. Kanatsoulis, X. Fu, N. D. Sidiropoulos, and M. Hong, "Structured SUMCOR multiview canonical analysis for large-scale data," *IEEE Transactions on Signal Processing*, vol. 67, no. 2, pp. 306–319, 2018.
- [19] A. Scaglione, R. Pagliari, and H. Krim, "The decentralized estimation of the sample covariance," in *2008 42nd Asilomar Conference on Signals, Systems and Computers*. IEEE, 2008, pp. 1722–1726.
- [20] L. Li, A. Scaglione, and J. H. Manton, "Distributed principal subspace estimation in wireless sensor networks," *IEEE Journal of Selected Topics in Signal Processing*, vol. 5, no. 4, pp. 725–738, 2011.
- [21] Y. Zeng and R. C. Hendriks, "Distributed delay and sum beamformer for speech enhancement via randomized gossip," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 22, no. 1, pp. 260–273, 2013.
- [22] X. Fu, K. Huang, E. E. Papalexakis, H. A. Song, P. Talukdar, N. D. Sidiropoulos, C. Faloutsos, and T. Mitchell, "Efficient and distributed generalized canonical correlation analysis for big multiview data," *IEEE Transactions on Knowledge and Data Engineering*, vol. 31, no. 12, pp. 2304–2318, 2018.
- [23] A. Bertrand and M. Moonen, "Distributed canonical correlation analysis in wireless sensor networks with application to distributed blind source separation," *IEEE Transactions on Signal Processing*, vol. 63, no. 18, pp. 4800–4813, 2015.
- [24] J. F. Mota, J. M. Xavier, P. M. Aguiar, and M. Puschel, "Distributed basis pursuit," *IEEE Transactions on Signal Processing*, vol. 60, no. 4, pp. 1942–1956, 2011.
- [25] J. Zhang, A. I. Koutrouvelis, R. Heusdens, and R. C. Hendriks, "Distributed rate-constrained LCMV beamforming," *IEEE Signal Processing Letters*, vol. 26, no. 5, pp. 675–679, 2019.
- [26] C. Hovine and A. Bertrand, "MAXVAR-based distributed correlation estimation in a wireless sensor network," *IEEE Transactions on Signal Processing*, vol. 70, pp. 5533–5548, 2022.
- [27] Y. Liu, X. Zhang, Y. Kang, L. Li, T. Chen, M. Hong, and Q. Yang, "Fedbcd: A communication-efficient collaborative learning framework for distributed features," *IEEE Transactions on Signal Processing*, vol. 70, pp. 4277–4290, 2022.
- [28] C. A. Musluoglu and A. Bertrand, "A unified algorithmic framework for distributed adaptive signal and feature fusion problems—part I: Algorithm derivation," *IEEE Transactions on Signal Processing*, 2023.
- [29] C. A. Musluoglu, C. Hovine, and A. Bertrand, "A unified algorithmic framework for distributed adaptive signal and feature fusion problems—part II: Convergence properties," *IEEE Transactions on Signal Processing*, 2023.
- [30] C. Hovine and A. Bertrand, "A distributed adaptive algorithm for non-smooth spatial filtering problems," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023, pp. 1–5.

- [31] —, “A distributed adaptive algorithm for non-smooth spatial filtering problems in wireless sensor networks,” *arXiv preprint arXiv:2403.08658*, 2024.
- [32] S. Kim, R. Pasupathy, and S. G. Henderson, “A guide to sample average approximation,” *Handbook of simulation optimization*, pp. 207–243, 2015.
- [33] S. V. Vaseghi, “Wiener filters,” *Advanced Signal Processing and Digital Noise Reduction*, pp. 140–163, 1996.
- [34] C. A. Musluoglu and A. Bertrand, “DASF toolbox,” https://github.com/AlexanderBertrandLab/DASF_toolbox, [Online].
- [35] F. H. Clarke, *Optimization and nonsmooth analysis*. SIAM, 1990.
- [36] R. T. Rockafellar and R. J.-B. Wets, *Variational analysis*. Springer Science & Business Media, 2009, vol. 317.
- [37] W. Karush, “Minima of functions of several variables with inequalities as side constraints,” *M. Sc. Dissertation. Dept. of Mathematics, Univ. of Chicago*, 1939.
- [38] H. Kuhn and A. Tucker, “Nonlinear programming,” in *Proc. 2nd Berkeley Symposium on Mathematical Statistics and Probabilistics*, 1951, pp. 481–492.
- [39] J. O. Royset and R. J. Wets, *An Optimization Primer*. Springer, 2021.
- [40] A. Beck, *First-order methods in optimization*. SIAM, 2017.
- [41] A. Themelis, “Proximal algorithms for structured nonconvex optimization,” 2018.
- [42] H. Attouch, J. Bolte, and B. F. Svaiter, “Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized Gauss–Seidel methods,” *Mathematical Programming*, vol. 137, no. 1, pp. 91–129, 2013.
- [43] J. Nocedal and S. J. Wright, *Numerical optimization*. Springer, 1999.
- [44] D. Charalambos and B. Aliprantis, *Infinite Dimensional Analysis: A Hitchhiker’s Guide*. Springer-Verlag Berlin and Heidelberg GmbH & Company KG, 2013.
- [45] W. Rudin *et al.*, *Principles of mathematical analysis*. McGraw-hill New York, 1976, vol. 3.
- [46] B. Ghojogh, F. Karray, and M. Crowley, “Eigenvalue and generalized eigenvalue problems: Tutorial,” *arXiv preprint arXiv:1903.11240*, 2019.
- [47] E. Anderson, Z. Bai, C. Bischof, L. S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney *et al.*, *LAPACK users’ guide*. SIAM, 1999.
- [48] W. Kozłowski, “The power method for the generalized eigenvalue problem,” *Mathematica Applicanda*, vol. 21, no. 35, 1992.
- [49] C. A. Musluoglu, M. Moonen, and A. Bertrand, “Improved tracking for distributed signal fusion optimization in a fully-connected wireless sensor network,” in *2022 30th European Signal Processing Conference (EUSIPCO)*. IEEE, 2022, pp. 1836–1840.
- [50] L. Armijo, “Minimization of functions having Lipschitz continuous first partial derivatives,” *Pacific Journal of mathematics*, vol. 16, no. 1, pp. 1–3, 1966.