



**HAL**  
open science

## Sparse phase retrieval via group-sparse optimization

Fabien Lauer, Henrik Ohlsson

► **To cite this version:**

Fabien Lauer, Henrik Ohlsson. Sparse phase retrieval via group-sparse optimization. 2014. hal-00951158

**HAL Id: hal-00951158**

**<https://hal.science/hal-00951158v1>**

Preprint submitted on 24 Feb 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Sparse phase retrieval via group-sparse optimization

F. Lauer<sup>1</sup> and H. Ohlsson<sup>2,3</sup>

<sup>1</sup> LORIA, Université de Lorraine, CNRS, Inria, France

<sup>2</sup> Dept. of Electrical Engineering and Computer Sciences, University of California, Berkeley, USA

<sup>3</sup> Dept. of Electrical Engineering, Linköping University, Sweden

February 24, 2014

## Abstract

This paper deals with sparse phase retrieval, i.e., the problem of estimating a vector from quadratic measurements under the assumption that few components are nonzero. In particular, we consider the problem of finding the sparsest vector consistent with the measurements and reformulate it as a group-sparse optimization problem with linear constraints. Then, we analyze the convex relaxation of the latter based on the minimization of a block  $\ell_1$ -norm and show various exact recovery and stability results in the real and complex cases. Invariance to circular shifts and reflections are also discussed for real vectors measured via complex matrices.

## 1 Introduction

The problem of recovering a signal from quadratic measurements is known as phase retrieval. A typical case with many applications, for instance in optics [23] or crystallography [11], is when the measurements correspond to the (squared) magnitude of the Fourier transform of the signal. Here, we consider the more general setting of arbitrary quadratic measurements,  $y_i = \mathbf{x}^H \mathbf{Q}_i \mathbf{x}$ ,  $i = 1, \dots, N$ , while focusing on the case where the signal  $\mathbf{x} \in \mathbb{C}^n$  is assumed to be sparse, i.e., with few nonzero entries in  $\mathbf{x}$ . As in compressive sensing [5, 1], which deals with the recovery of sparse signals from linear measurements, the sparsity prior reduces the number of measurements required to recover the signal.

**Related work.** Seminal works on phase retrieval [13, 9, 10, 8] did not consider the sparsity prior. Though these methods were able to incorporate prior information on the support, they were not designed to estimate the support with limited information on its size. More recently, matrix lifting techniques were developed in [3, 2, 22] for phase retrieval and more particularly for the case of sparse signals in [21, 12, 18]. The basic idea is to apply a change of variable resulting in linearized measurement equations with a rank-1 constraint on the new matrix variable  $\mathbf{X} = \mathbf{x}\mathbf{x}^H$ . Then, the rank-1 constraint is relaxed to the problem of minimizing the rank of the matrix, which is further relaxed to the minimization of the nuclear norm. In these methods, the sparsity prior is typically incorporated as the minimization of an  $\ell_1$ -norm, which induces a trade-off between the satisfaction of the rank-1 constraint and the sparsity of the solution. Other methods focusing on sparsity are typically iterative, like Fienup-type methods [15] using alternate projections or the GESPAR method [20] implementing a local search strategy with a bi-directional greedy algorithm. These iterative methods usually come without recovery guarantees.

Finally, note that sparse phase retrieval also enters the more general framework of nonlinear compressed sensing, as investigated in [17] for analytic functions computing the measurements, in [7] for quasi-linear functions and in [14] for polynomials.

**Contribution.** We propose a convex approach to the sparse phase retrieval problem. This approach relies on two main steps: the linearization of the constraints inducing a group-sparse structure on the variables and a convex relaxation of the group-sparse optimization problem enforcing this structure. More precisely, the linearization is based on the Veronese map lifting the signal to a higher dimensional space. This map is invariant to a global sign change and preserves sparsity in the sense that the lifted signal has a group-sparsity similar to the sparsity of the original signal. Thus, the proposed method amounts to estimating a group-sparse signal satisfying linear constraints, from which the original signal can be recovered. This estimation relies on a convex relaxation of the group-sparse problem based on a sum of norms, or block  $\ell_1$ -norm. Thus, while methods based on matrix lifting lead to semi-definite programming problems, the proposed approach yields a more amenable second-order cone programming formulation. This formulation is also easily extended to deal with noisy measurements. In addition to these algorithmic benefits, we derive exact recovery conditions in the noiseless case and stable recovery guarantees in the presence of noise.

Note that the approach taken here is similar in spirit to the one derived in [14] for the more general problem of finding sparse solutions of polynomial systems of equations. However, the analysis in [14] is limited to the real case and does not apply to polynomials without linear terms as the ones found in phase retrieval.

**Paper organization.** For the sake of clarity, we first detail in Sect. 2 the proposed method in the real case before extending it in Sect. 3 to the complex case. The effect of noise and stability results are discussed in Sect. 4. The case of real signals measured via complex vectors is considered in Sect. 5 which also deals with the invariance of the measurements to circular shifts and reflections. Finally, Section 6 tests the proposed methods in numerical experiments.

**Notations.** Matrices are written with bold uppercase letters and vectors in bold lowercase letters, except for the  $i$ th column  $\mathbf{A}_i$  of a matrix  $\mathbf{A}$ . The notation  $(\mathbf{A})_{i,j}$  denotes the element at the  $i$ th row and  $j$ th column of a matrix  $\mathbf{A}$ .  $\Re(\cdot)$  and  $\Im(\cdot)$  denote the real and imaginary parts of a complex number, vector or matrix, and  $i$  the imaginary unit. The superscripts  $T$  and  $H$  denote the transpose and conjugate transpose, respectively, i.e.,  $\mathbf{z}^H = \bar{\mathbf{z}}^T$ .  $\|\cdot\|_p$  denotes the  $\ell_p$ -norm in  $\mathbb{R}^n$ , while  $\|\cdot\|$  denotes the norm in  $\mathbb{C}^n$  induced by the inner product as  $\|\mathbf{z}\| = \sqrt{\mathbf{z}^H \mathbf{z}}$ . The  $\ell_0$ -pseudo-norm of a real or complex vector  $\mathbf{x}$  of dimension  $n$  is defined as  $\|\mathbf{x}\|_0 = |\{j \in \{1, \dots, n\} : x_j \neq 0\}|$  and denotes the number of nonzero components  $x_j$ . We also define the  $\ell_0$ -pseudo-norm of a vector-valued sequence  $\{\mathbf{u}_i\}_{i=1}^N$  as  $\|\{\mathbf{u}_i\}_{i=1}^N\|_0 = |\{i \in \{1, \dots, N\} : \mathbf{u}_i \neq \mathbf{0}\}|$ .

## 2 The real case

We write the sparse phase retrieval problem as

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{x}\|_0 \\ \text{s.t. } y_i = (\mathbf{q}_i^T \mathbf{x})^2 = \mathbf{x}^T \mathbf{Q}_i \mathbf{x}, \quad i = 1, \dots, N, \end{aligned} \tag{1}$$

where  $y_i \in \mathbb{R}$  are the measurements and  $\mathbf{Q}_i = \mathbf{q}_i \mathbf{q}_i^T \in \mathbb{R}^{n \times n}$ . Due to symmetry, solutions to (1) are defined up to their sign and the goal is to obtain an estimate  $\hat{\mathbf{x}} = \pm \mathbf{x}_0$ , for  $\mathbf{x}_0$  in the solution set of (1). In particular, we are interested in the case where (1) has a unique pair of solutions  $\{\mathbf{x}_0, -\mathbf{x}_0\}$ , while conditions ensuring such a uniqueness are discussed in [16, 19].

The proposed method relies on two subsequent relaxations. While the first one linearizes the constraints, the second one convexifies the objective function.

### 2.1 First level of relaxation

Let the Veronese map of degree 2,  $\nu : \mathbb{R}^n \rightarrow \mathbb{R}^M$ , be defined by

$$\nu(\mathbf{x}) = [x_1^2, x_1 x_2, \dots, x_2^2, x_2 x_3, \dots, x_{n-1}^2, x_{n-1} x_n, x_n^2]^T,$$

and the subscript  $i_j$  denote the index of its component equal to  $x_i x_j$ , i.e.,

$$i_j = \sum_{k=1}^{\min\{i,j\}-1} (n-k+1) + |j-i| + 1 = \sum_{k=1}^{\min\{i,j\}-1} (n-k) + \min\{i,j\} + |j-i|. \quad (2)$$

This notation is symmetric, i.e.,  $i_j$  and  $j_i$  denote the same index, and will be used throughout the paper to index the components of vectors of  $\mathbb{R}^M$  or  $\mathbb{C}^M$ .

The constraints of the phase retrieval problem (1) can be rewritten as

$$\mathbf{A}\nu(\mathbf{x}) = \mathbf{y}$$

with  $\mathbf{A} \in \mathbb{R}^{N \times M}$  and  $M = \binom{n+1}{2} = n(n+1)/2$ .

Let  $\mathbf{W}_j$  be an  $n \times M$ -binary matrix such that<sup>1</sup>  $\mathbf{W}_j \nu(\mathbf{x}) = x_j \mathbf{x}$ , i.e.,  $\mathbf{W}_j \nu(\mathbf{x})$  is the vector of  $n$  entries corresponding to the monomials in  $\nu(\mathbf{x})$  including  $x_j$ . Then, we have

$$x_j = 0 \Leftrightarrow \mathbf{W}_j \nu(\mathbf{x}) = \mathbf{0}$$

and the objective function in (1) can be written as

$$\|\mathbf{x}\|_0 = \|\{\mathbf{W}_j \nu(\mathbf{x})\}_{j=1}^n\|_0. \quad (3)$$

Thus, (1) can be reformulated as the nonlinear group-sparse optimization problem

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \quad & \|\{\mathbf{W}_j \nu(\mathbf{x})\}_{j=1}^n\|_0 \\ \text{s.t.} \quad & \mathbf{A}\nu(\mathbf{x}) = \mathbf{y}. \end{aligned} \quad (4)$$

Note that  $\nu(\mathbf{x}) = \nu(-\mathbf{x})$ , but that for  $\mathbf{x} \neq \pm \mathbf{x}_0$ ,  $\nu(\mathbf{x}) \neq \nu(\mathbf{x}_0)$ . Thus, the problem can be posed as the one of recovering the value of  $\nu(\mathbf{x}_0)$ , from which  $\mathbf{x}_0$  can be inferred up to its sign.

To estimate  $\nu(\mathbf{x}_0)$ , we relax (4) to

$$\begin{aligned} \min_{\mathbf{v} \in \mathbb{R}^M} \quad & \|\{\mathbf{W}_j \mathbf{v}\}_{j=1}^n\|_0 \\ \text{s.t.} \quad & \mathbf{A}\mathbf{v} = \mathbf{y} \\ & v_{jj} \geq 0, \quad j = 1, \dots, n, \end{aligned} \quad (5)$$

where the variables in  $\mathbf{v}$  estimating the components of  $\nu(\mathbf{x})$  are not constrained to be interdependent monomials of  $n$  base variables, but the last constraints in (5) nonetheless ensure that the  $v_{jj}$ 's estimating the  $x_j^2$ 's are positive.

## 2.2 Convex relaxation

Problem (5) is a (linear) group-sparse optimization problem with highly overlapping groups. While groupwise-greedy algorithms, such as the one proposed in [14], can be applied, their analysis is not available for the case of overlapping groups. Therefore, here, we consider the convex relaxation approach which aims at solving (5) via the following surrogate formulation:

$$\begin{aligned} \hat{\mathbf{v}} = \arg \min_{\mathbf{v} \in \mathbb{R}^M} \quad & \sum_{j=1}^n \|\mathbf{W}_j \mathbf{W} \mathbf{v}\|_2 \\ \text{s.t.} \quad & \mathbf{A}\mathbf{v} = \mathbf{y} \\ & v_{jj} \geq 0, \quad j = 1, \dots, n, \end{aligned} \quad (6)$$

where we introduced the diagonal matrix  $\mathbf{W}$  of precompensating weights  $(\mathbf{W})_{i,i} = w_i = \|\mathbf{A}_i\|_2$ , and which can be solved efficiently by off-the-shelf Second-Order Cone Programming (SOCP) solvers.

<sup>1</sup>More precisely, the entries of  $\mathbf{W}_j$  are given by  $(\mathbf{W}_j)_{k,l} = \delta_{l,jk}$ ,  $l = 1, \dots, M$ ,  $k = 1, \dots, n$ , where  $\delta$  is the Kronecker delta and  $j_k$  is an index as in (2).

Then, we easily obtain an estimate of  $\mathbf{x}_0$  from the estimate  $\hat{\mathbf{v}}$  of  $\nu(\mathbf{x})$  as  $\hat{\mathbf{x}} = \nu^{-1}(\hat{\mathbf{v}})$ , where the inverse mapping  $\nu^{-1}$  is defined as

$$\nu^{-1}(\mathbf{v}) = \begin{cases} \frac{1}{\sqrt{v_{ii}}} [v_{1i}, v_{2i}, \dots, v_{ni}]^T, & \text{if } i > 0 \text{ and } \frac{v_{ji}^2}{v_{ii}} = v_{jj}, \forall j \in \{1, \dots, n\} \\ \mathbf{0}, & \text{otherwise,} \end{cases}$$

where

$$i = \begin{cases} \min_{j \in \{1, \dots, n\}} j, \text{ s.t. } v_{jj} > 0, \text{ if } \exists j \text{ such that } v_{jj} > 0 \\ 0, & \text{otherwise.} \end{cases}$$

In the above, the first nonzero entry (of index  $i$ ) is assumed to be positive to fix the signs and make  $\nu^{-1}$  injective. This definition also ensures that

$$\nu^{-1}(\nu(\mathbf{x})) = \pm \mathbf{x}$$

and that

$$\|\nu^{-1}(\mathbf{v})\|_0 \leq \|\{\mathbf{W}_j \mathbf{v}\}_{j=1}^n\|_0 \quad (7)$$

since, for  $\nu^{-1}(\mathbf{v}) \neq \mathbf{0}$ ,  $\mathbf{W}_j \mathbf{v} = \mathbf{0} \Rightarrow v_{ji} = 0 \Rightarrow (\nu^{-1}(\mathbf{v}))_j = 0$ .

### 2.3 Analysis

We now turn to theoretical guarantees offered by the proposed approach. First, the following theorem provides the rationale for tackling the sparse phase retrieval problem via the group-sparse optimization formulation (5).

**Theorem 1.** *If the solution  $\mathbf{v}^*$  to (5) is unique and yields  $\mathbf{x}^* = \nu^{-1}(\mathbf{v}^*) \neq \mathbf{0}$  such that  $y_i = (\mathbf{x}^*)^T \mathbf{Q}_i \mathbf{x}^*$ ,  $i = 1, \dots, N$ , then  $\{\mathbf{x}^*, -\mathbf{x}^*\}$  is the unique pair of solutions of (1).*

*Proof.* Assume there is an  $\mathbf{x}_0 \neq \pm \mathbf{x}^*$  satisfying the constraints of (1) and at least as sparse as  $\mathbf{x}^*$ . Then,  $\mathbf{A}\nu(\mathbf{x}_0) = \mathbf{y}$  and, by using (3) and (7),

$$\|\{\mathbf{W}_j \nu(\mathbf{x}_0)\}_{j=1}^n\|_0 = \|\mathbf{x}_0\|_0 \leq \|\mathbf{x}^*\|_0 \leq \|\{\mathbf{W}_j \mathbf{v}^*\}_{j=1}^n\|_0,$$

which contradicts the fact that  $\mathbf{v}^*$  is the unique solution to (4) unless  $\nu(\mathbf{x}_0) = \mathbf{v}^*$ . But since  $\mathbf{x}_0 \neq \pm \mathbf{x}^*$ , we have  $x_{0j}^2 \neq (x_j^*)^2$  for some  $j \in \{1, \dots, n\}$ , which implies  $(\nu(\mathbf{x}_0))_{jj} \neq (\nu(\mathbf{x}^*))_{jj} = (\nu(\nu^{-1}(\mathbf{v}^*)))_{jj}$ . Therefore, by using Lemma 1 in Appendix A with the assumption  $\mathbf{x}^* = \nu^{-1}(\mathbf{v}^*) \neq \mathbf{0}$ , there cannot be such an  $\mathbf{x}_0$ .  $\square$

The following results regarding the convex formulation (6) are based on the notion of mutual coherence.

**Definition 1.** *The mutual coherence of a matrix  $\mathbf{A} = [\mathbf{A}_1, \dots, \mathbf{A}_M] \in \mathbb{R}^{N \times M}$  is*

$$\mu(\mathbf{A}) = \max_{1 \leq i < j \leq M} \frac{|\mathbf{A}_i^T \mathbf{A}_j|}{\|\mathbf{A}_i\|_2 \|\mathbf{A}_j\|_2}.$$

With this definition, we can state an exact recovery result (proof given in Appendix B.1).

**Theorem 2.** *Let  $\mathbf{x}_0$  be such that  $y_i = \mathbf{x}_0^T \mathbf{Q}_i \mathbf{x}_0$ ,  $i = 1, \dots, N$ , and  $\mathbf{v}_0 = \nu(\mathbf{x}_0)$ . If the condition*

$$\|\mathbf{x}_0\|_0 < \frac{1}{2\sqrt{n}} \sqrt{1 + \frac{1}{\mu^2(\mathbf{A})}}$$

*holds, then  $\mathbf{v}_0$  is the unique solution to (6).*

**Corollary 1.** Let  $\mathbf{x}_0$  be a feasible point of (1). If the condition

$$\|\mathbf{x}_0\|_0 < \frac{1}{2\sqrt{n}} \sqrt{1 + \frac{1}{\mu^2(\mathbf{A})}}$$

holds, then  $\{\mathbf{x}_0, -\mathbf{x}_0\}$  is the unique pair of solutions to the minimization problem (1) and they can be computed as  $\mathbf{x}_0 = \pm\nu^{-1}(\hat{\mathbf{v}})$  with  $\hat{\mathbf{v}}$  the solution to (6).

*Proof.* Assume there exists another solution  $\mathbf{x}_1 \neq \pm\mathbf{x}_0$  to (1), and thus with  $\|\mathbf{x}_1\|_0 \leq \|\mathbf{x}_0\|_0$ . Then, Theorem 2 implies that both  $\nu(\mathbf{x}_1)$  and  $\nu(\mathbf{x}_0)$  are *unique* solutions to (6) and thus that  $\nu(\mathbf{x}_1) = \nu(\mathbf{x}_0) = \hat{\mathbf{v}}$ . But this contradicts the definition of the mapping  $\nu$  implying  $\nu(\mathbf{x}_1) \neq \nu(\mathbf{x}_0)$  whenever  $\mathbf{x}_1 \neq \pm\mathbf{x}_0$ . Therefore the assumption  $\mathbf{x}_1 \neq \pm\mathbf{x}_0$  cannot hold and  $\{\mathbf{x}_0, -\mathbf{x}_0\}$  is the unique pair of solutions to (1), while  $\nu^{-1}(\hat{\mathbf{v}}) = \nu^{-1}(\nu(\mathbf{x}_0)) = \pm\mathbf{x}_0$ .  $\square$

### 3 The complex case

Consider now the problem in complex domain:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{C}^n} \|\mathbf{x}\|_0 \\ \text{s.t. } y_i = |\mathbf{q}_i^H \mathbf{x}|^2, \quad i = 1, \dots, N, \end{aligned} \quad (8)$$

where  $y_i \in \mathbb{R}$  and  $\mathbf{q}_i \in \mathbb{C}^n$ .

The equations in the problem above are invariant to multiplication by a unit complex scalar  $z$  with  $|z| = 1$ . Thus, there are sets of solutions of the form  $T(\mathbf{x}_0) = \{\mathbf{x} \in \mathbb{C}^n : \mathbf{x} = z\mathbf{x}_0, z \in \mathbb{C}, |z| = 1\}$ , and the goal is to obtain an estimate  $\hat{\mathbf{x}} \in T(\mathbf{x}_0)$ , from which  $T(\mathbf{x}_0) = T(\hat{\mathbf{x}})$  can be inferred due to the property of the invariance set:

$$\forall \mathbf{x} \in T(\mathbf{x}_0), \quad T(\mathbf{x}) = T(\mathbf{x}_0), \quad (9)$$

which can be proved as follows. Let  $\mathbf{x} = z\mathbf{x}_0$  with  $|z| = 1$ , then  $T(\mathbf{x}) = \{\mathbf{a} \in \mathbb{C}^n : \mathbf{a} = b\mathbf{x}, b \in \mathbb{C}, |b| = 1\} = \{\mathbf{a} \in \mathbb{C}^n : \mathbf{a} = bz\mathbf{x}_0, b \in \mathbb{C}, |b| = 1\} = \{\mathbf{a} \in \mathbb{C}^n : \mathbf{a} = c\mathbf{x}_0, c \in \mathbb{C}, |c| = 1\} = T(\mathbf{x}_0)$ .

#### 3.1 First level of relaxation

As for the real case, the linearization of the equations will use the Veronese map, which we redefine for complex vectors as follows.

**Definition 2** (Complex Veronese map). *The complex Veronese map  $\nu : \mathbb{C}^n \rightarrow \mathbb{C}^M$  is defined by*

$$\nu(\mathbf{x}) = [x_1\bar{x}_1, x_1\bar{x}_2, \dots, x_2\bar{x}_2, x_2\bar{x}_3, \dots, x_{n-1}\bar{x}_{n-1}, x_{n-1}\bar{x}_n, x_n\bar{x}_n]^T,$$

for which the subscript  $ij$ , defined as in (2), denotes the component index such that  $(\nu(\mathbf{x}))_{ij}$  equals either  $x_i\bar{x}_j$  or  $x_j\bar{x}_i$  (note that, for all pairs  $(i, j)$ , there is exactly one such component).

Note that  $\nu(\mathbf{x}) = \nu(\mathbf{x}')$  for all  $\mathbf{x}' \in T(\mathbf{x})$ , since for  $1 \leq i \leq j \leq n$ ,  $x'_i\bar{x}'_j = zx_j\bar{z}x_i = |z|^2 x_i\bar{x}_j = x_i\bar{x}_j$ , but that  $\mathbf{x}' \notin T(\mathbf{x})$  implies  $\nu(\mathbf{x}) \neq \nu(\mathbf{x}')$ , since  $\mathbf{x}' \notin T(\mathbf{x}) \Rightarrow |x'_j| \neq |x_j| \Rightarrow x'_j\bar{x}'_j \neq x_j\bar{x}_j$ .

Then, we define the inverse mapping as follows.

**Definition 3** (Inverse complex Veronese map). *The inverse complex Veronese map,  $\nu^{-1} : \mathbb{C}^M \rightarrow \mathbb{C}^n$ , is defined by*

$$\nu^{-1}(\mathbf{v}) = \begin{cases} \frac{1}{\sqrt{v_{ii}}} [v_{1i}, v_{2i}, \dots, v_{ni}]^T, & \text{if } i > 0 \text{ and } \frac{|v_{ji}|^2}{v_{ii}} = v_{jj}, \forall j \in \{1, \dots, n\} \\ \mathbf{0}, & \text{otherwise,} \end{cases}$$

where

$$i = \begin{cases} \min_{j \in \{1, \dots, n\}} j, \text{ s.t. } \Re(v_{jj}) > 0, \Im(v_{jj}) = 0, \text{ if } \exists j \text{ such that } \Re(v_{jj}) > 0, \Im(v_{jj}) = 0 \\ 0, & \text{otherwise.} \end{cases}$$

In particular, we have  $\nu^{-1}(\nu(\mathbf{x})) \in T(\mathbf{x})$ .

Note that the square root acts on a real and positive number  $v_{ii}$  and is thus also a real positive number. This implies  $x_i = \sqrt{v_{ii}} \in \mathbb{R}^+$  and that we arbitrarily set  $\Im(x_i) = 0$  to fix the value of  $z$  in the equation  $\nu^{-1}(\nu(\mathbf{x})) = z\mathbf{x}$ , for a complex number  $z$  with  $|z| = 1$ , and thus make  $\nu^{-1}$  injective.

With these definitions at hand, the equations in (8) are reformulated via Lemma 3 (all Lemmas are given in Appendix A) as follows:

$$y_i = \mathbf{x}^H \mathbf{q}_i \mathbf{q}_i^H \mathbf{x} = 2\Re(\nu(\mathbf{q}_i)^H \nu(\mathbf{x})) - \sum_{j=1}^n (\nu(\mathbf{q}_i))_{jj} (\nu(\mathbf{x}))_{jj}, \quad i = 1, \dots, N.$$

Define the vectors  $\mathbf{a}_i \in \mathbb{C}^M$ ,  $i = 1, \dots, N$ , with components given by  $(\mathbf{a}_i)_{jk} = 2(\nu(\mathbf{q}_i))_{jk}$  for  $1 \leq j < k \leq n$  and  $(\mathbf{a}_i)_{jj} = (\nu(\mathbf{q}_i))_{jj}$ ,  $j = 1, \dots, n$ . Then, the equations above, linear wrt. to  $\nu(\mathbf{x})$ , can be rewritten as  $y_i = \Re(\mathbf{a}_i^H \nu(\mathbf{x}))$ ,  $i = 1, \dots, N$ .

Additionally define the binary matrices  $\mathbf{W}_j$  such that the vector  $\mathbf{W}_j \nu(\mathbf{x})$  contains all the monomials of  $\nu(\mathbf{x})$  including either  $x_j$  or  $\bar{x}_j$ :

$$\mathbf{W}_j \nu(\mathbf{x}) = [x_1 \bar{x}_j, x_2 \bar{x}_j, \dots, x_{j-1} \bar{x}_j, x_j \bar{x}_j, x_j \bar{x}_{j+1}, \dots, x_j \bar{x}_n]^T \in \mathbb{C}^n.$$

Then, we have

$$\|\mathbf{x}\|_0 = \|\{\mathbf{W}_j \nu(\mathbf{x})\}_{j=1}^n\|_0 \quad (10)$$

and problem (8) can be rewritten as the nonlinear group-sparse optimization program

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{C}^n} \quad & \|\{\mathbf{W}_j \nu(\mathbf{x})\}_{j=1}^n\|_0 \\ \text{s.t.} \quad & \mathbf{y} = \Re(\mathbf{A}\nu(\mathbf{x})), \end{aligned} \quad (11)$$

where  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_N]^H$ .

Next, we relax this formulation by substituting  $\mathbf{v} \in \mathbb{C}^M$  for  $\nu(\mathbf{x})$ , which yields

$$\begin{aligned} \min_{\mathbf{v} \in \mathbb{C}^M} \quad & \|\{\mathbf{W}_j \mathbf{v}\}_{j=1}^n\|_0 \\ \text{s.t.} \quad & \mathbf{y} = \Re(\mathbf{A}\mathbf{v}) \\ & v_{jj} \in \mathbb{R}^+, \quad j = 1, \dots, n, \end{aligned} \quad (12)$$

where the last constraints ensure that the  $v_{jj}$ 's estimating the modulus of the base variables are positive real numbers.

The following theorem shows that this relaxation can be used as a proxy to solve the original problem.

**Theorem 3.** *If the solution  $\mathbf{v}^*$  to (12) is unique and yields  $\mathbf{x}^* = \nu^{-1}(\mathbf{v}^*) \neq \mathbf{0}$  such that  $y_i = |\mathbf{q}_i^H \mathbf{x}^*|^2$ ,  $i = 1, \dots, N$ , then  $T(\mathbf{x}^*)$  is the unique set of solutions of (8).*

*Proof.* Assume there is an  $\mathbf{x}_0 \neq T(\mathbf{x}^*)$  satisfying the constraints of (8) and at least as sparse as  $\mathbf{x}^*$ . Then  $\mathbf{y} = \Re(\mathbf{A}\nu(\mathbf{x}_0))$ , and, by using (10) and Lemma 4,

$$\|\{\mathbf{W}_j \nu(\mathbf{x}_0)\}_{j=1}^n\|_0 = \|\mathbf{x}_0\|_0 \leq \|\mathbf{x}^*\|_0 \leq \|\{\mathbf{W}_j \mathbf{v}^*\}_{j=1}^n\|_0,$$

which contradicts the fact that  $\mathbf{v}^*$  is the unique solution to (12) unless  $\nu(\mathbf{x}_0) = \mathbf{v}^*$ . But since  $\mathbf{x}_0 \notin T(\mathbf{x}^*)$ , we have  $|x_{0j}| \neq |x_j^*|$  and thus  $x_{0j} \bar{x}_{0j} \neq x_j^* \bar{x}_j^*$  for some  $j \in \{1, \dots, n\}$ , which implies  $(\nu(\mathbf{x}_0))_{jj} \neq (\nu(\mathbf{x}^*))_{jj} = (\nu(\nu^{-1}(\mathbf{v}^*)))_{jj}$ . Therefore, by using Lemma 5 with the assumption  $\mathbf{x}^* = \nu^{-1}(\mathbf{v}^*) \neq \mathbf{0}$ , there cannot be such an  $\mathbf{x}_0$ .  $\square$

### 3.2 Convex relaxation

We now introduce a second level of relaxation by replacing the  $\ell_0$ -pseudo norm by a block- $\ell_1$  norm. This leads to a convex relaxation in the form of a SOCP:

$$\begin{aligned} \min_{\mathbf{v} \in \mathbb{C}^M} \quad & \sum_{j=1}^n \|\mathbf{W}_j^R \Re(\mathbf{v}) + i\mathbf{W}_j^I \Im(\mathbf{v})\| \\ \text{s.t.} \quad & \mathbf{y} = \Re(\mathbf{A}\mathbf{v}) \\ & v_{jj} \in \mathbb{R}^+, \quad j = 1, \dots, n, \end{aligned} \quad (13)$$

where  $\mathbf{W}_j^R = \mathbf{W}_j \mathbf{W}^R$  and  $\mathbf{W}_j^I = \mathbf{W}_j \mathbf{W}^I$  with  $\tilde{\mathbf{W}} = \begin{pmatrix} \mathbf{W}^R & \mathbf{0} \\ \mathbf{0} & \mathbf{W}^I \end{pmatrix}$  a diagonal matrix of precompensating weights given by  $(\mathbf{W}^R)_{i,i} = \|\Re(\mathbf{A}_i)\|_2$  and  $(\mathbf{W}^I)_{i,i} = \|\Im(\mathbf{A}_i)\|_2$ .

**Theorem 4.** Let  $\mathbf{x}_0$  be such that  $y_i = \mathbf{x}_0^H \mathbf{q}_i \mathbf{q}_i^H \mathbf{x}_0$ ,  $i = 1, \dots, N$ ,  $\mathbf{v}_0 = \nu(\mathbf{x}_0)$  and  $\tilde{\mathbf{A}} = [\Re(\mathbf{A}), -\Im(\mathbf{A})]$ . If the condition

$$\|\mathbf{x}_0\|_0 < \frac{1}{2\sqrt{2n}} \sqrt{1 + \frac{1}{\mu^2(\tilde{\mathbf{A}})}}$$

holds, then  $\mathbf{v}_0$  is the unique solution to (13).

*Proof.* The vector  $\mathbf{v}_0$  is the unique solution to (13) if the inequality

$$\sum_{j=1}^n \|\mathbf{W}_j^R \Re(\mathbf{v}_0 + \boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\mathbf{v}_0 + \boldsymbol{\delta})\| > \sum_{j=1}^n \|\mathbf{W}_j^R \Re(\mathbf{v}_0) + i\mathbf{W}_j^I \Im(\mathbf{v}_0)\|$$

holds for all  $\boldsymbol{\delta} \in \mathbb{C}^M$  such that  $\Re(\mathbf{A}(\mathbf{v}_0 + \boldsymbol{\delta})) = \mathbf{y}$ , which implies the constraint  $\Re(\mathbf{A}\boldsymbol{\delta}) = \mathbf{0}$  on  $\boldsymbol{\delta}$ . The inequality above can be rewritten as

$$\sum_{j \in I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\| + \sum_{j \notin I_0} \|\mathbf{W}_j^R \Re(\mathbf{v}_0 + \boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\mathbf{v}_0 + \boldsymbol{\delta})\| - \|\mathbf{W}_j^R \Re(\mathbf{v}_0) + i\mathbf{W}_j^I \Im(\mathbf{v}_0)\| > 0,$$

where  $I_0 = \{j \in \{1, \dots, n\} : \mathbf{W}_j^R \Re(\mathbf{v}_0) = \mathbf{0} \wedge \mathbf{W}_j^I \Im(\mathbf{v}_0) = \mathbf{0}\}$ . By the triangle inequality,  $\|\mathbf{a} + \mathbf{b}\| - \|\mathbf{a}\| \geq -\|\mathbf{b}\|$  with  $\mathbf{a} = \mathbf{W}_j^R \Re(\mathbf{v}_0) + i\mathbf{W}_j^I \Im(\mathbf{v}_0)$ , this condition is met if

$$\sum_{j \in I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\| - \sum_{j \notin I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\| > 0$$

or

$$\sum_{j=1}^n \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\| - 2 \sum_{j \notin I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\| > 0. \quad (14)$$

By defining  $G_j$  as the set of indexes corresponding to nonzero columns of  $\mathbf{W}_j$ , Lemma 6 yields

$$\begin{aligned} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|^2 &= \sum_{i \in G_j} (w_i^R)^2 \Re(\delta_i)^2 + (w_i^I)^2 \Im(\delta_i)^2 \\ &\leq 2n \frac{\mu^2(\tilde{\mathbf{A}})}{1 + \mu^2(\tilde{\mathbf{A}})} \|\mathbf{W}^R \Re(\boldsymbol{\delta}) + i\mathbf{W}^I \Im(\boldsymbol{\delta})\|^2. \end{aligned}$$

Due to the fact that  $\bigcup_{k \in \{1, \dots, n\}} G_k = \{1, \dots, M\}$ , we also have

$$\begin{aligned} \|\mathbf{W}^R \Re(\boldsymbol{\delta}) + i\mathbf{W}^I \Im(\boldsymbol{\delta})\|^2 &= \sum_{i=1}^M (w_i^R)^2 \Re(\delta_i)^2 + (w_i^I)^2 \Im(\delta_i)^2 \\ &\leq \sum_{k=1}^n \sum_{i \in G_k} (w_i^R)^2 \Re(\delta_i)^2 + (w_i^I)^2 \Im(\delta_i)^2 \\ &= \sum_{k=1}^n \|\mathbf{W}_k^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_k^I \Im(\boldsymbol{\delta})\|^2 \\ &\leq \left( \sum_{k=1}^n \|\mathbf{W}_k^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_k^I \Im(\boldsymbol{\delta})\| \right)^2, \end{aligned}$$

which then leads to

$$\|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|^2 \leq 2n \frac{\mu^2(\tilde{\mathbf{A}})}{1 + \mu^2(\tilde{\mathbf{A}})} \left( \sum_{k=1}^n \|\mathbf{W}_k^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_k^I \Im(\boldsymbol{\delta})\| \right)^2.$$



Introducing this result in (14) gives the condition

$$\sum_{j=1}^n \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\| - 2(n - |I_0|) \frac{\mu(\tilde{\mathbf{A}})\sqrt{2n}}{\sqrt{1 + \mu^2(\tilde{\mathbf{A}})}} \sum_{k=1}^n \|\mathbf{W}_k^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_k^I \Im(\boldsymbol{\delta})\| > 0.$$

Finally, given that  $|I_0| = n - \|\mathbf{x}_0\|_0$ , this yields

$$\sum_{j=1}^n \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\| - 2\|\mathbf{x}_0\|_0 \frac{\mu(\tilde{\mathbf{A}})\sqrt{2n}}{\sqrt{1 + \mu^2(\tilde{\mathbf{A}})}} \sum_{k=1}^n \|\mathbf{W}_k^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_k^I \Im(\boldsymbol{\delta})\| > 0. \quad (15)$$

or, for  $\boldsymbol{\delta} \neq \mathbf{0}$ ,

$$\|\mathbf{x}_0\|_0 < \frac{\sqrt{1 + \mu^2(\tilde{\mathbf{A}})}}{2\mu(\tilde{\mathbf{A}})\sqrt{2n}},$$

which can be rewritten as in the statement of the Theorem.  $\square$

**Corollary 2.** *Let  $\mathbf{x}_0$  be a feasible point of (8). If the condition*

$$\|\mathbf{x}_0\|_0 < \frac{1}{2\sqrt{2n}} \sqrt{1 + \frac{1}{\mu^2(\tilde{\mathbf{A}})}}$$

*holds, then  $T(\mathbf{x}_0)$  is the unique set of solutions to the minimization problem (8) and it can be computed as  $T(\mathbf{x}_0) = T(\nu^{-1}(\hat{\mathbf{v}}))$  with  $\hat{\mathbf{v}}$  the solution to (13).*

*Proof.* Assume there exists another solution  $\mathbf{x}_1 \notin T(\mathbf{x}_0)$  to (8), and thus with  $\|\mathbf{x}_1\|_0 \leq \|\mathbf{x}_0\|_0$ . Then, Theorem 4 implies that both  $\nu(\mathbf{x}_1)$  and  $\nu(\mathbf{x}_0)$  are *unique* solutions to (13) and thus that  $\nu(\mathbf{x}_1) = \nu(\mathbf{x}_0) = \hat{\mathbf{v}}$ . But this contradicts Definition 2 implying  $\nu(\mathbf{x}_1) \neq \nu(\mathbf{x}_0)$  whenever  $\mathbf{x}_1 \notin T(\mathbf{x}_0)$ . Therefore the assumption  $\mathbf{x}_1 \notin T(\mathbf{x}_0)$  cannot hold and  $T(\mathbf{x}_0)$  is the unique set of solutions to (8), while  $\nu^{-1}(\hat{\mathbf{v}}) = \nu^{-1}(\nu(\mathbf{x}_0)) \in T(\mathbf{x}_0)$ , which, by using (9), implies  $T(\mathbf{x}_0) = T(\nu^{-1}(\hat{\mathbf{v}}))$ .  $\square$

## 4 Stable recovery in the presence of noise

Consider now the case where the measurements  $\mathbf{y}$  are perturbed by an additive noise  $\mathbf{e} \in \mathbb{R}^N$  of bounded  $\ell_2$ -norm,  $\|\mathbf{e}\|_2 \leq \varepsilon$ . Then, the equations in (1) and (8) are of the form  $y_i = |\mathbf{q}_i^H \mathbf{x}|^2 + e_i$  with the noise terms  $e_i$  to be estimated together with the sparse signal. Note that in this context, multiple solutions with different noise vectors can be valid. Thus, we aim at stability results bounding the error on the estimates by a function of  $\varepsilon$  rather exact recovery ones. Details on the proposed method to achieve these goals are given below, first for real signals and then for complex ones.

### 4.1 Stability in the real case

In the noisy case with real data, the problem that we need to solve becomes

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{e} \in \mathbb{R}^N} \quad & \|\mathbf{x}\|_0 \\ \text{s.t.} \quad & y_i = |\mathbf{q}_i^T \mathbf{x}|^2 + e_i, \quad i = 1, \dots, N, \\ & \|\mathbf{e}\|_2 \leq \varepsilon, \end{aligned} \quad (16)$$

where  $y_i \in \mathbb{R}$  and  $\mathbf{q}_i \in \mathbb{R}^n$ . Following the approach of Sect. 2 leads to a convex relaxation in the form of the SOCP

$$\begin{aligned} \min_{\mathbf{v} \in \mathbb{R}^M} \quad & \sum_{j=1}^n \|\mathbf{W}_j \mathbf{W} \mathbf{v}\|_2 \\ \text{s.t.} \quad & \|\mathbf{y} - \mathbf{A} \mathbf{v}\|_2 \leq \varepsilon \\ & v_{jj} \geq 0, \quad j = 1, \dots, n, \end{aligned} \quad (17)$$

for which we have the following stability result.

**Theorem 5.** Let  $(\mathbf{x}_0, \mathbf{e}_0)$  denote a solution to (16). If the inequality

$$\|\mathbf{x}_0\|_0 < \frac{1}{2n^2(n+1)} \left(1 + \frac{1}{\mu(\mathbf{A})}\right) \quad (18)$$

holds, then the solution  $\hat{\mathbf{v}}$  to (17) must obey

$$\|\mathbf{W}(\hat{\mathbf{v}} - \nu(\mathbf{x}_0))\|_2^2 \leq \frac{4n\varepsilon^2}{1 - \mu(\mathbf{A})[2n^2(n+1)\|\mathbf{x}_0\|_0 - 1]}. \quad (19)$$

If, in addition,  $\varepsilon = 0$ , then  $\hat{\mathbf{x}} = \pm\mathbf{x}_0$ .

We omit the proof which is similar to the one of Theorem 6 in [14], except for the last statement concluding on the stability of  $\hat{\mathbf{x}}$ , and which closely follows the one for the complex case of Theorem 6 to be detailed below.

## 4.2 Stability in the complex case

Consider now the complex variant of the problem perturbed by noise:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{C}^n, \mathbf{e} \in \mathbb{R}^N} \quad & \|\mathbf{x}\|_0 \\ \text{s.t.} \quad & y_i = |\mathbf{q}_i^H \mathbf{x}|^2 + e_i, \quad i = 1, \dots, N, \\ & \|\mathbf{e}\|_2 \leq \varepsilon, \end{aligned} \quad (20)$$

where  $y_i \in \mathbb{R}$  and  $\mathbf{q}_i \in \mathbb{C}^n$ . The solution to this problem can be approached via the convex relaxation

$$\begin{aligned} \min_{\mathbf{v} \in \mathbb{C}^M} \quad & \sum_{j=1}^n \|\mathbf{W}_j^R \Re(\mathbf{v}) + \mathbf{iW}_j^I \Im(\mathbf{v})\| \\ \text{s.t.} \quad & \|\mathbf{y} - \Re(\mathbf{A}\mathbf{v})\|_2 \leq \varepsilon \\ & v_{jj} \in \mathbb{R}^+, \quad j = 1, \dots, n. \end{aligned} \quad (21)$$

As for the real case, we have a stability result for the estimation of  $\nu(\mathbf{x})$ .

**Theorem 6.** Let  $(\mathbf{x}_0, \mathbf{e}_0)$  denote a solution to (20). If the inequality

$$\|\mathbf{x}_0\|_0 < \frac{1}{2n^2(n+1)} \left(1 + \frac{1}{\mu(\tilde{\mathbf{A}})}\right) \quad (22)$$

holds, then the solution  $\hat{\mathbf{v}}$  to (21) must obey

$$\|\mathbf{W}^R \Re(\hat{\mathbf{v}} - \nu(\mathbf{x}_0)) + \mathbf{iW}^I \Im(\hat{\mathbf{v}} - \nu(\mathbf{x}_0))\|_2^2 \leq \frac{4n\varepsilon^2}{1 - \mu(\tilde{\mathbf{A}})[2n^2(n+1)\|\mathbf{x}_0\|_0 - 1]}. \quad (23)$$

If, in addition,  $\varepsilon = 0$ , then  $\hat{\mathbf{x}} = \nu^{-1}(\hat{\mathbf{v}}) \in T(\mathbf{x}_0)$ .

*Proof.* Assume (20) has a solution  $(\mathbf{x}_0, \mathbf{e}_0)$ . Let define  $\mathbf{v}_0 = \nu(\mathbf{x}_0)$  and  $\boldsymbol{\delta} = \hat{\mathbf{v}} - \mathbf{v}_0$ . The proof follows a path similar to that of Theorem 3.1 in [6], which was adapted in [14] to the group-sparse setting and which is here further extended to the complex case.

Due to the definition of  $\hat{\mathbf{v}}$  as a minimizer of (21),  $\boldsymbol{\delta}$  must satisfy either

$$\sum_{j=1}^n \|\mathbf{W}_j^R \Re(\mathbf{v}_0 + \boldsymbol{\delta}) + \mathbf{iW}_j^I \Im(\mathbf{v}_0 + \boldsymbol{\delta})\| < \sum_{j=1}^n \|\mathbf{W}_j^R \Re(\mathbf{v}_0) + \mathbf{iW}_j^I \Im(\mathbf{v}_0)\|$$

or  $\boldsymbol{\delta} = \mathbf{0}$ , in which case the statement is obvious. The inequality above can be rewritten as

$$\sum_{j \in I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + \mathbf{iW}_j^I \Im(\boldsymbol{\delta})\| + \sum_{j \notin I_0} \|\mathbf{W}_j^R \Re(\mathbf{v}_0 + \boldsymbol{\delta}) + \mathbf{iW}_j^I \Im(\mathbf{v}_0 + \boldsymbol{\delta})\| - \|\mathbf{W}_j^R \Re(\mathbf{v}_0) + \mathbf{iW}_j^I \Im(\mathbf{v}_0)\| < 0,$$

where  $I_0 = \{j \in \{1, \dots, n\} : \mathbf{W}_j^R \Re(\mathbf{v}_0) = \mathbf{0} \wedge \mathbf{W}_j^I \Im(\mathbf{v}_0) = \mathbf{0}\}$ . By the triangle inequality,  $\|\mathbf{a} + \mathbf{b}\| - \|\mathbf{a}\| \geq -\|\mathbf{b}\|$ , this implies

$$\sum_{j \in I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\| - \sum_{j \notin I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\| < 0. \quad (24)$$

In addition,  $\boldsymbol{\delta}$  must satisfy the constraints in (21) as

$$\|\Re(\mathbf{A}(\mathbf{v}_0 + \boldsymbol{\delta})) - \mathbf{y}\|_2 \leq \varepsilon,$$

in which  $\mathbf{y}$  can be replaced by  $\Re(\mathbf{A}\mathbf{v}_0) + \mathbf{e}_0$ , leading to

$$\|\Re(\mathbf{A}\boldsymbol{\delta}) - \mathbf{e}_0\|_2 \leq \varepsilon.$$

Using  $\|\mathbf{a}\|_2 \leq \|\mathbf{a} - \mathbf{b}\|_2 + \|\mathbf{b}\|_2$ , this implies  $\|\Re(\mathbf{A}\boldsymbol{\delta})\|_2 \leq 2\varepsilon$ , which further gives

$$\begin{aligned} 4\varepsilon^2 &\geq \|\Re(\mathbf{A}\boldsymbol{\delta})\|_2^2 = \|\tilde{\mathbf{A}}\tilde{\boldsymbol{\delta}}\|_2^2 = \|\tilde{\mathbf{A}}\tilde{\mathbf{W}}^{-1}\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2 = (\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}})^T \tilde{\mathbf{W}}^{-1} \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} \tilde{\mathbf{W}}^{-1} (\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}) \\ &= \|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2 + (\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}})^T (\tilde{\mathbf{W}}^{-1} \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} \tilde{\mathbf{W}}^{-1} - \mathbf{I}) (\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}) \\ &\geq \|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2 - \left| (\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}})^T (\tilde{\mathbf{W}}^{-1} \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} \tilde{\mathbf{W}}^{-1} - \mathbf{I}) (\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}) \right| \\ &\geq \|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2 - |\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}|^T |\tilde{\mathbf{W}}^{-1} \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} \tilde{\mathbf{W}}^{-1} - \mathbf{I}| |\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}| \\ &\geq \|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2 - \mu(\tilde{\mathbf{A}}) (\|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_1^2 - \|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2) \\ &= (1 + \mu(\tilde{\mathbf{A}})) \|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2 - \mu(\tilde{\mathbf{A}}) \|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_1^2 \end{aligned} \quad (25)$$

where we used  $\tilde{\mathbf{W}}^{-1} \tilde{\mathbf{W}} = \mathbf{I}$  and the fact that the diagonal entries of  $|\tilde{\mathbf{W}}^{-1} \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} \tilde{\mathbf{W}}^{-1} - \mathbf{I}|$  are zeros while off-diagonal entries are bounded from above by  $\mu(\tilde{\mathbf{A}})$ .

Due to  $\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})$  being a vector with a subset of entries from  $\mathbf{W}^R \Re(\boldsymbol{\delta}) + i\mathbf{W}^I \Im(\boldsymbol{\delta})$ , we have  $\|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2 = \|\mathbf{W}^R \Re(\boldsymbol{\delta}) + i\mathbf{W}^I \Im(\boldsymbol{\delta})\|^2 \geq \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|^2$ ,  $j = 1, \dots, n$ , and thus

$$\|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2 \geq \frac{1}{n} \sum_{j=1}^n \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|^2. \quad (26)$$

Since the groups defined by the  $\mathbf{W}_j$ 's overlap,  $\|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2 = \|\mathbf{W}^R \Re(\boldsymbol{\delta}) + i\mathbf{W}^I \Im(\boldsymbol{\delta})\| \leq \sum_{j=1}^n \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|$ , and the squared  $\ell_1$ -norm in (25) can be bounded by

$$\|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_1^2 \leq M \|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2 \leq M \left( \sum_{j=1}^n \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\| \right)^2. \quad (27)$$

Introducing the bounds (26)–(27) in (25) yields

$$\frac{1 + \mu(\tilde{\mathbf{A}})}{n} \sum_{j=1}^n \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|^2 - \mu(\tilde{\mathbf{A}}) M \left( \sum_{j=1}^n \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\| \right)^2 \leq 4\varepsilon^2. \quad (28)$$

We will now use this inequality to derive an upper bound on  $\sum_{j=1}^n \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|^2$ , which will also apply to  $\|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2 \leq \sum_{j=1}^n \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|^2 = \sum_{j=1}^n (\|\mathbf{W}_j^R \Re(\boldsymbol{\delta})\|_2^2 + \|\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|_2^2)$ , since the groups overlap and the squared components of  $\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}$  are summed multiple times in the right-hand side. To derive the upper bound, we first introduce a few notations:

$$a = \sum_{j \in I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|, \quad b = \sum_{j \notin I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|,$$

and

$$c_0 = \frac{\sum_{j \in I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|^2}{\left( \sum_{j \in I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\| \right)^2} \in \left[ \frac{1}{|I_0|}, 1 \right],$$

$$c_1 = \frac{\sum_{j \notin I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|^2}{\left(\sum_{j \notin I_0} \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|\right)^2} \in \left[\frac{1}{n - |I_0|}, 1\right],$$

where the box bounds are obtained by classical relations between the  $\ell_1$  and  $\ell_2$  norms<sup>2</sup>. With these notations, the term to bound is rewritten as

$$\sum_{j=1}^n \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|^2 = c_0 a^2 + c_1 b^2,$$

while the inequality (28) becomes

$$\frac{1 + \mu(\tilde{\mathbf{A}})}{n} (c_0 a^2 + c_1 b^2) - \mu(\tilde{\mathbf{A}}) M (a + b)^2 \leq 4\varepsilon^2.$$

We further reformulate this constraint by letting  $a = \rho b$ :

$$\frac{1 + \mu(\tilde{\mathbf{A}})}{n} (c_0 \rho^2 + c_1) b^2 - \mu(\tilde{\mathbf{A}}) M (1 + \rho)^2 b^2 \leq 4\varepsilon^2. \quad (29)$$

Let  $\gamma = (1 + \rho)^2 / (c_0 \rho^2 + c_1)$ . Due to (24), we have  $a < b$  and thus  $\rho \in [0, 1)$ , which, together with the bounds on  $c_0$  and  $c_1$ , gives the constraints  $1 \leq \gamma \leq 4(n - |I_0|)$ . By setting  $V = (c_0 \rho^2 + c_1) b^2$ , (29) is rewritten as

$$\frac{1 + \mu(\tilde{\mathbf{A}})}{n} V - \mu(\tilde{\mathbf{A}}) M \gamma V \leq 4\varepsilon^2,$$

where

$$\frac{1 + \mu(\tilde{\mathbf{A}})}{n} - \mu(\tilde{\mathbf{A}}) M \gamma \geq \frac{1 + \mu(\tilde{\mathbf{A}})}{n} - 4(n - |I_0|) \mu(\tilde{\mathbf{A}}) M > 0,$$

since  $\gamma \leq 4(n - |I_0|)$  and the positivity is ensured by the condition (22) and the fact that  $\|\mathbf{x}_0\|_0 = n - |I_0|$ . Thus,

$$\|\tilde{\mathbf{W}} \tilde{\boldsymbol{\delta}}\|_2^2 \leq \sum_{j=1}^n \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|^2 = V \leq \frac{4n\varepsilon^2}{1 + \mu(\tilde{\mathbf{A}}) - 4\mu(\tilde{\mathbf{A}})nM\|\mathbf{x}_0\|_0},$$

which proves the stability result in (23) since  $\|\tilde{\mathbf{W}} \tilde{\boldsymbol{\delta}}\|_2^2 = \|\mathbf{W}^R \Re(\boldsymbol{\delta}) + i\mathbf{W}^I \Im(\boldsymbol{\delta})\|^2$ .

To conclude in the case  $\varepsilon = 0$ , it remains to see that (23) implies  $\hat{\mathbf{v}} = \nu(\mathbf{x}_0)$  and that Definition 3 ensures  $\nu^{-1}(\nu(\mathbf{x}_0)) \in T(\mathbf{x}_0)$ .  $\square$

## 5 Complex data, but real solutions

Consider now the following problem:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{x}\|_0 \\ \text{s.t. } y_i = |\mathbf{q}_i^H \mathbf{x}|^2, \quad i = 1, \dots, N, \end{aligned} \quad (30)$$

where the signal  $\mathbf{x} \in \mathbb{R}^n$  and measurements  $y_i \in \mathbb{R}$  are assumed to be real while the vectors  $\mathbf{q}_i \in \mathbb{C}^n$  can be complex. In this case,  $\mathbf{v} = \nu(\mathbf{x})$  is also a real vector and solutions can be approximated via a dedicated version of (13):

$$\begin{aligned} \min_{\mathbf{v} \in \mathbb{R}^M} \sum_{j=1}^n \|\mathbf{W}_j^R \mathbf{v}\|_2 \\ \text{s.t. } \mathbf{y} = \Re(\mathbf{A} \mathbf{v}) = \Re(\mathbf{A}) \mathbf{v} \\ v_{jj} \geq 0, \quad j = 1, \dots, n, \end{aligned} \quad (31)$$

<sup>2</sup>Let  $\mathbf{u} \in \mathbb{R}^{|I_0|}$  with  $u_j = \|\mathbf{W}_j^R \Re(\boldsymbol{\delta}) + i\mathbf{W}_j^I \Im(\boldsymbol{\delta})\|$ . Then,  $c_0 = (\|\mathbf{u}\|_2 / \|\mathbf{u}\|_1)^2$  and the bounds are obtained by the classical relation  $\forall \mathbf{u} \in \mathbb{R}^k, \|\mathbf{u}\|_2 \leq \|\mathbf{u}\|_1 \leq \sqrt{k} \|\mathbf{u}\|_2$ .

where  $\mathbf{A} \in \mathbb{C}^{N \times M}$  is defined as in Sect. 3.1 and  $\mathbf{W}_j^R = \mathbf{W}_j \mathbf{W}^R$  with precompensating weights given by  $(\mathbf{W}^R)_{i,i} = \|\Re(\mathbf{A}_i)\|_2$ .

For this particular case, Theorem 7 below is similar in spirit to Theorem 4, but allows us to gain a  $\sqrt{2}$  factor by using Lemma 7 instead of Lemma 6 in order to take into account that  $\boldsymbol{\delta}$  belongs to  $\mathbb{R}^M$  (detailed proof given in Appendix B.2). This results in a bound on  $\|\mathbf{x}_0\|_0$  similar to the one in Theorem 2 for the real case and based on the mutual coherence of the real part of  $\mathbf{A}$ . Since  $\mu(\Re(\mathbf{A})) \leq \mu(\tilde{\mathbf{A}})$ , this also improves (relaxes) the bound compared with the one of Theorem 4.

**Theorem 7.** *Let  $\mathbf{x}_0 \in \mathbb{R}^n$  be such that  $y_i = \mathbf{x}_0^T \mathbf{q}_i \mathbf{q}_i^H \mathbf{x}_0$ ,  $i = 1, \dots, N$ , and  $\mathbf{v}_0 = \nu(\mathbf{x}_0)$ . If the condition*

$$\|\mathbf{x}_0\|_0 < \frac{1}{2\sqrt{n}} \sqrt{1 + \frac{1}{\mu^2(\Re(\mathbf{A}))}}$$

holds, then  $\mathbf{v}_0$  is the unique solution to (31).

In a typical instance of Problem (30), the measurements  $\mathbf{y}$  correspond to the power spectrum of a real signal. However, in this case, Theorem 7 does not apply since the solution is known not to be unique due to invariances of the measurements. The following subsections first describe a practical technique to deal with such invariances and then focus on the power spectrum case.

## 5.1 Invariances

Consider the case where the measurements are invariant to some transformation of the signal. A typical example is when  $\mathbf{x} \in \mathbb{R}^n$  and  $\{\mathbf{q}_i \in \mathbb{C}^n\}$  forms a Fourier basis. Then, the equations  $y_i = |\mathbf{q}_i^H \mathbf{x}|^2$  are invariant to circular shifts of the components of  $\mathbf{x}$ . This is problematic since shifted versions of  $\mathbf{x}$  lead to shuffled<sup>3</sup> versions of  $\nu(\mathbf{x})$ . Thus, multiple shuffled  $\mathbf{v}$ 's satisfy the linearized constraints,  $y_i = \Re(\mathbf{A})\mathbf{v}$ , and so does any convex combination of them. These convex combinations need not be sparse as they combine vectors with different sparsity patterns, but lead to lower or equal values of the convex cost function of (31).

To circumvent this issue, we need to linearize the constraints by a shift-invariant transformation  $\phi$ , such that  $\phi(\mathbf{x}_0) = \phi(\mathbf{x}_1)$  for  $\mathbf{x}_0$  and  $\mathbf{x}_1$  two shifted versions of the same vector. To lead to an effective estimation method, the transformation  $\phi$  must also retain sparsity in a sense similar to  $\nu$ .

Consider the transformation  $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^M$  defined by<sup>4</sup>

$$\phi(\mathbf{x}) = \nu(\text{shift}(\mathbf{x}, 1 - k)), \quad \text{with } k = \arg \max_{i \in \{1, \dots, n\}} |x_i|, \quad (32)$$

where  $\text{shift}(\cdot, k)$  stands for the  $k$ -steps circular shift operator. The transformation  $\phi$  in (32) first shifts the vector  $\mathbf{x}$  such that the first component is the one with maximal magnitude. This results in a shift-invariant transformation which retains sparsity and the linearization ability via the mapping  $\nu$ . The linearized constraints remain the same, i.e.,  $y_i = \Re(\mathbf{A})\phi$ , but another constraint must be added to (31) in order to account for the shift-invariance.

More precisely, the definition of  $\phi$  ensures that

$$(\phi(\mathbf{x}))_{11} \geq (\phi(\mathbf{x}))_{jj}, \quad j = 2, \dots, n.$$

Thus, a valid vector  $\hat{\phi}$  estimating  $\phi(\mathbf{x}_0)$  can be found by solving

$$\begin{aligned} \hat{\phi} &= \arg \min_{\phi \in \mathbb{R}^M} \sum_{j=2}^n \|\mathbf{W}_j^R \phi\|_2 \\ \text{s.t. } &\mathbf{y} = \Re(\mathbf{A}) \phi \\ &\phi_{11} \geq \phi_{jj} \geq 0, \quad j = 2, \dots, n. \end{aligned} \quad (33)$$

<sup>3</sup>Circular shifts of  $\mathbf{x}$  lead to rearrangements of the components of  $\nu(\mathbf{x})$  which are not exactly circular shifts. For example, with  $\mathbf{x}_0 = [1, 1, 0, 0]^T$  and  $\mathbf{x}_1 = [0, 1, 1, 0]^T$ , we have  $\nu(\mathbf{x}_0) = [1, 1, 0, 0, 1, 0, 0, 0, 0]^T$  and  $\nu(\mathbf{x}_1) = [0, 0, 0, 0, 1, 1, 0, 1, 0, 0]^T$ .

<sup>4</sup>In the case where the argmax in (32) is not a singleton,  $k$  is arbitrarily set to the minimum of the indexes in the argmax and  $\phi$  cannot be shift-invariant. Assumptions regarding this issue will be made clear in Proposition 1 and Theorem 8 below.

Note that the the cost function does not involve the first group of variables since  $x_1$  is assumed to be nonzero. In comparison with (31), this formulation is still convex but cannot have multiple solutions that are shifted/shuffled versions of one another.

Finally, the set of shifted solutions  $\{\mathbf{x}_k\}_{k=1}^n$  to (30) is approximated by the set

$$\phi^{-1}(\hat{\phi}) = \left\{ \text{shift}(\nu^{-1}(\hat{\phi}), k) \right\}_{k=1}^n.$$

Following a similar approach, we can additionally take into account reflections by defining

$$\begin{cases} \phi(\mathbf{x}) = \nu(\mathbf{x}_2) \\ \mathbf{x}_2 = \varphi(\mathbf{x}) = \begin{cases} \mathbf{x}_1, & \text{if } \sum_{i=2}^{n/2} |x_{1i}|^2 \geq \sum_{i=2+\frac{n}{2}}^n |x_{1i}|^2 \\ \text{shift}(\text{reflection}(\mathbf{x}_1), 1), & \text{otherwise} \end{cases} \\ \mathbf{x}_1 = \text{shift}(\mathbf{x}, 1 - k) \\ k = \arg \max_{i \in \{1, \dots, n\}} |x_i|, \end{cases} \quad (34)$$

where we assume  $n$  to be even and  $\text{reflection}(\cdot)$  is the reflection operator defined by  $\text{reflection}(\mathbf{x}) = [x_n, x_{n-1}, \dots, x_2, x_1]^T$ . In plain words, the transformation  $\phi$  in (34) first shifts the vector  $\mathbf{x}$  such that the first component is the one with maximal magnitude. Then, it applies a centered reflection to the shifted  $\mathbf{x}$ , named  $\mathbf{x}_1$ , only if the sum of squares over the first entries of  $\mathbf{x}_1$  (without the first one) is smaller than the one over the last ones. If this is the case, the result is shifted again to recover the first component of  $\mathbf{x}_2$  with maximal magnitude. Finally, the Veronese map is applied to  $\mathbf{x}_2$  to give  $\phi(\mathbf{x})$ .

Invariance of  $\phi = \nu \circ \varphi$  to circular shifts and reflections is implied by the invariance of  $\varphi$  given in the Proposition below (proof in Appendix B.3).

**Proposition 1.** *For all  $\mathbf{x} \in \mathbb{C}^n$  such that  $|\arg \max_{i \in \{1, \dots, n\}} |x_i|| = 1$ , the following statements hold for the transformation  $\varphi$  defined in (34):*

1.  $\varphi$  is idempotent, i.e.,  $\varphi \circ \varphi(\mathbf{x}) = \varphi(\mathbf{x})$ ;

and, for all  $\mathbf{x}$  additionally satisfying  $\sum_{i=2}^{n/2} |x_{1i}|^2 \neq \sum_{i=2+\frac{n}{2}}^n |x_{1i}|^2$  with  $\mathbf{x}_1 = \text{shift}(\mathbf{x}, \arg \max_{i \in \{1, \dots, n\}} |x_i|)$ ,

2.  $\varphi$  is shift-invariant, i.e.,  $\forall s \in \mathbb{Z}$ ,  $\varphi(\mathbf{x}) = \varphi(\text{shift}(\mathbf{x}, s))$ ;

3.  $\varphi$  is reflection-invariant, i.e.,  $\varphi(\mathbf{x}) = \varphi(\text{reflection}(\mathbf{x}))$ .

Note that statements 2 and 3 imply that  $\varphi$  is invariant to any combination of shifts and reflections.

Proposition 1 also shows the idempotence of  $\varphi$ , which of course does not transfer to  $\phi$  directly but which is however very useful. Indeed, this allows us to test if a candidate  $\mathbf{x}_2$ , supposed to be the result of  $\varphi$  applied to some vector, is consistent with the definition of  $\varphi$  as  $\mathbf{x}_2 \stackrel{?}{=} \varphi(\mathbf{x}_2)$ , which can be checked via simple inequalities. Since these inequalities only involve the (squared) magnitude of the entries in the vector, they can also be easily embedded as linear constraints in (33) to compute an estimate  $\hat{\phi}$  that is consistent with the transformation (34). This yields the convex program

$$\begin{aligned} \hat{\phi} = \arg \min_{\phi \in \mathbb{R}^M} \sum_{j=2}^n \|\mathbf{W}_j^R \phi\|_2 & \quad (35) \\ \text{s.t. } \mathbf{y} = \mathfrak{R}(\mathbf{A}) \phi & \quad (\text{data fitting}) \\ \phi_{11} \geq \phi_{jj} \geq 0, \quad j = 2, \dots, n & \quad (\text{shift-invariance}) \\ \sum_{i=2}^{n/2} \phi_{ii} \geq \sum_{i=2+\frac{n}{2}}^n \phi_{ii} & \quad (\text{reflection-invariance}), \end{aligned}$$

where  $\phi_{jj}$  estimates  $|x_j|^2$ . The advantage of using (35) is that it has a single solution independently of the number of shifted/reflected solutions to (30), in the sense of the next theorem.

**Theorem 8.** *If Problem (30) has a unique set of shifted/reflected solutions  $S(\mathbf{x}_0) = \{\text{shift}(\mathbf{x}_0, k) : k \in \{1, \dots, n\}\} \cup \{\text{shift}(\text{reflection}(\mathbf{x}_0), k) : k \in \{1, \dots, n\}\}$ , and if  $|\arg \max_{i \in \{1, \dots, n\}} |x_{0i}|| = 1$ , then there is exactly one vector  $\nu(\mathbf{x})$  with  $\mathbf{x} \in S(\mathbf{x}_0)$  in the feasible set of (35).*

*Proof.* Since all  $\mathbf{x} \in S(\mathbf{x}_0)$  are solutions to (30), they satisfy the first constraint in (35) as  $\mathbf{y} = \Re(\mathbf{A})\nu(\mathbf{x})$ . By Proposition 1,  $\varphi$  is constant over the set  $S(\mathbf{x}_0)$  and equal to one of the vectors from this set, say  $\mathbf{x}_0$ , for which we have in particular  $\mathbf{x}_0 = \varphi(\mathbf{x}_0)$ . Then,  $\nu(\mathbf{x}_0) = \nu \circ \varphi(\mathbf{x}_0) = \phi(\mathbf{x}_0)$  and thus  $\nu(\mathbf{x}_0)$  satisfies all the constraints of (35) and is a feasible point.

Now take an  $\mathbf{x} \in S(\mathbf{x}_0)$  with  $\mathbf{x} \neq \mathbf{x}_0$ . Then, either  $\mathbf{x} = \text{shift}(\mathbf{x}_0, k)$  or  $\mathbf{x} = \text{shift}(\text{reflection}(\mathbf{x}_0), k)$  for some  $k \in \{1, \dots, n\}$ . This implies  $\arg \max_i (\nu(\mathbf{x}))_{ii} = \arg \max_i |x_i| \neq \arg \max_i |x_{0i}| = \arg \max_i (\nu(\mathbf{x}_0))_{ii}$ . Since  $\mathbf{x}_0$  satisfies the constraints in (35), we have  $\arg \max_i (\nu(\mathbf{x}_0))_{ii} = 1$  and thus  $\arg \max_i (\nu(\mathbf{x}))_{ii} \neq 1$ , which shows that  $\nu(\mathbf{x})$  violates the second constraint and is not a feasible point of (35).  $\square$

**Illustrative example.** Consider the vectors ( $\mathbf{x}_1$  and  $\mathbf{x}_2$  are not related to the notations of (34))

$$\begin{aligned} \mathbf{x}_1 &= [1, 2, 3, 4, 0, 0]^T, & \mathbf{x}_2 &= [4, 0, 0, 1, 2, 3]^T, & \mathbf{x}_3 &= [0, 0, 4, 3, 2, 1]^T, \\ \mathbf{x}_4 &= [2, 1, 0, 0, 4, 3]^T, & \mathbf{x}_5 &= [4, 3, 2, 1, 0, 0]^T, \end{aligned}$$

which are all obtained by shifts and reflections of the same vector. They all lead to the same  $\phi(\mathbf{x}_i) = \nu(\mathbf{x}_5)$ . Indeed, the transformation  $\phi$  first applies the required shift and reflect operations to map  $\mathbf{x}_i$  to  $\mathbf{x}_5$  which has a first component with the largest magnitude and the largest half-sum over its first entries. Then,  $\phi$  computes the Veronese map of  $\mathbf{x}_5$ . Also note that all the  $\nu(\mathbf{x}_i)$  are feasible with respect to  $\mathbf{y} = \Re(\mathbf{A})\nu(\mathbf{x}_i)$ , but only  $\nu(\mathbf{x}_5)$  satisfies the additional constraints implementing shift and reflection invariance in (35). Therefore, we can obtain  $\hat{\phi} = \nu(\mathbf{x}_5)$  as the unique solution to (35).

## 5.2 Support recovery from the power spectrum

When the measurements correspond to the squared magnitude of the Fourier transform of the signal, i.e., its (squared) power spectrum, there is another issue beside shift/reflection-invariances. Even with the correct support  $\mathcal{S} = \text{supp}(\mathbf{x}_0)$ , the linear system  $\mathbf{y} = \Re(\mathbf{A}_{\mathcal{S}})\phi_{\mathcal{S}}$ , where  $\mathbf{A}_{\mathcal{S}}$  is the submatrix of  $\mathbf{A}$  with the corresponding columns, is under-determined. More precisely, for all  $i$  and all  $j$ ,  $|q_{ij}| = 1$ , which implies  $(\mathbf{a}_i)_{jj} = (\nu(\mathbf{q}_i))_{jj} = 1$  and  $\mathbf{A}_{jj} = \mathbf{1}$ . Thus, even when limited to the support  $\mathcal{S}$ ,  $\mathbf{A}_{\mathcal{S}}$  has  $|\mathcal{S}|$  similar columns and is rank-deficient.

Therefore, in this case the proposed approach cannot exactly recover  $\mathbf{x}_0$  and (35) is used to estimate the support  $\mathcal{S}$ . Indeed, knowing the correct support can be useful for other methods dedicated to the classical phase retrieval problem [8].

Though this case seems unfavorable, if the number of measurements satisfies  $N \geq 2n - 1$ , the autocorrelation of  $\mathbf{x}_0$  (with zero padding),

$$r_k = \sum_{i=1}^n x_i x_{i+k}, \quad k = -n + 1, \dots, n - 1,$$

can be computed via the inverse Fourier transform of  $\mathbf{y}$  and thus is available. This information can be included as linear constraints on the components  $\phi_{i(i+k)}$  of  $\phi$  estimating the cross products  $x_i x_{i+k}$  to drive the solution towards a satisfactory one.

In addition, following [12, 20], this can be used to restrict the support of the solution as follows. First, note that with  $N \geq 2n - 1$ , the measurements are not invariant to circular shifts of  $\mathbf{x}$  but of  $\mathbf{x}$  with zero-padding. Thus, we cannot assume  $|x_1| \geq |x_j|$ ,  $j = 2, \dots, n$ . However, we can fix  $x_1 \neq 0$ . Then, for all base variable index  $j \in \{1, \dots, n\}$ ,

$$r_k = 0, \quad k = j - 1, \dots, n - 1 \quad \Rightarrow \quad x_j = 0. \quad (36)$$

Thus, the corresponding variables in  $\phi$  can be set to zero, i.e.,  $\mathbf{W}_j \phi = \mathbf{0}$ . Moreover, let  $j_m$  be the minimum of the  $j$ 's satisfying the condition in (36), then  $x_{j_m-1} \neq 0$ , which can be favored

by removing the corresponding term in the cost function. Combining all the information on the solution in a convex program leads to the final formulation

$$\begin{aligned}
\hat{\phi} = \arg \min_{\phi \in \mathbb{R}^M} & \sum_{j=2}^{j_m-2} \|\mathbf{W}_j^R \phi\|_2 & (37) \\
\text{s.t. } & \mathbf{y} = \mathfrak{R}(\mathbf{A}) \phi & \text{(data fitting)} \\
& \phi_{jj} \geq 0, \quad j = 1, \dots, j_m - 1 & \text{(structural knowledge)} \\
& \sum_{i=2}^{n/2} \phi_{ii} \geq \sum_{i=2+\frac{n}{2}}^n \phi_{ii} & \text{(reflection-invariance)} \\
& \sum_{i=1}^{n-k} \phi_{i(i+k)} = r_k, \quad k = 0, \dots, n-1 & \text{(autocorrelation)} \\
& \mathbf{W}_j \phi = \mathbf{0}, \quad j = j_m, \dots, n & \text{(restricted support)}.
\end{aligned}$$

A noise-tolerant version of (37) is obtained by replacing the constraint  $\mathbf{y} = \mathfrak{R}(\mathbf{A}) \phi$  by  $\|\mathbf{y} - \mathfrak{R}(\mathbf{A}) \phi\|_2 \leq \varepsilon$ . This modification also applies to (31) and (35).

## 6 Experiments

For the experiments, we consider two variants of the proposed approach: the convex method described in details in the previous sections and the greedy method for solving the group-sparse optimization problems (5) and (12). Implementation details for these two methods are as described in [14], with slight modifications to handle complex variables. In particular, the convex method uses the iterative reweighting of [4] adapted to the group-sparse setting to enhance the sparsity of the solution. The greedy method starts with an empty support and, at each iteration, adds the group of variables in  $\mathbf{v}$  corresponding to the base variable  $x_j$  that results in the best approximation of  $\mathbf{y}$ .

For complex signals,  $\mathbf{x}_0$ , we measure the relative error corresponding to the normalized distance between the estimate  $\hat{\mathbf{x}}$  and the set  $T(\mathbf{x}_0)$ , i.e.,  $\min_{\mathbf{x} \in T(\mathbf{x}_0)} \|\hat{\mathbf{x}} - \mathbf{x}\| / \|\mathbf{x}\|$ . Exact recovery is detected when this error is smaller than  $10^{-6} \|\mathbf{x}_0\|$ .

**Exact recovery in the noiseless case.** We estimate the probability of exact recovery at various sparsity levels,  $\|\mathbf{x}_0\|_0$ , in the complex case where  $n = 20$  and  $N = 50$ . For each sparsity level, the probability is estimated as the percentage of successful trials over a Monte Carlo experiment with 100 trials. In each trial,  $N$  complex vectors  $\mathbf{q}_i \in \mathbb{C}^n$  are drawn from a zero-mean Gaussian distribution of unit variance to measure a random signal  $\mathbf{x}_0 \in \mathbb{C}^n$  with  $\|\mathbf{x}_0\|_0$  nonzero entries at random locations whose values are drawn from a zero-mean Gaussian distribution of unit variance. Results shown in Fig. 1 indicate that the proposed approach can exactly recover sufficiently sparse signals with high probability.

Similar experiments are performed to evaluate the influence of the number of measurements  $N$  on the probability of exact recovery. Results reported in the right plot of Fig. 1 show that, with a sparsity level of  $\|\mathbf{x}_0\|_0 = 4$ , the convex method requires more measurements for perfect recovery than the greedy strategy. However,  $N \approx \|\mathbf{x}_0\|_0^3$  measurements are already sufficient to exactly recover  $\mathbf{x}_0$  in all trials.

**Stable recovery in the noisy case.** We now consider the noisy case where  $y_i = |\mathbf{q}_i^H \mathbf{x}_0|^2 + e_i$ ,  $i = 1, \dots, N$ , and  $\|e\|_2 \leq \varepsilon$ . Over 100 trials, Figure 2 reports the mean relative error and the rate of successful support recovery for  $N = 50$  and  $\varepsilon = 3$ . These results show that sufficiently sparse signals can be accurately estimated from noisy quadratic measurements and that the probability of support recovery in this case follows a curve similar to the one of the probability of exact recovery in the noiseless case. This means that the methods are robust to noise regarding the recovery of the correct support.



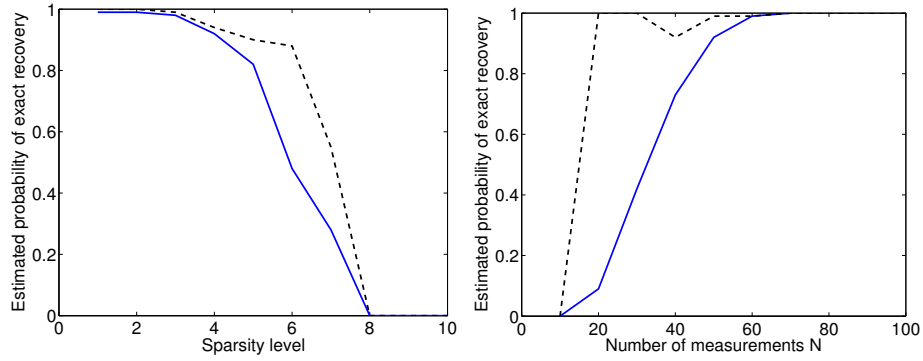


Figure 1: Estimated probability of exact recovery in the noiseless case for the convex relaxation (13) (plain line) and the greedy method applied to (12) (dashed line) versus the number of nonzeros (left) and the number of measurements (right).

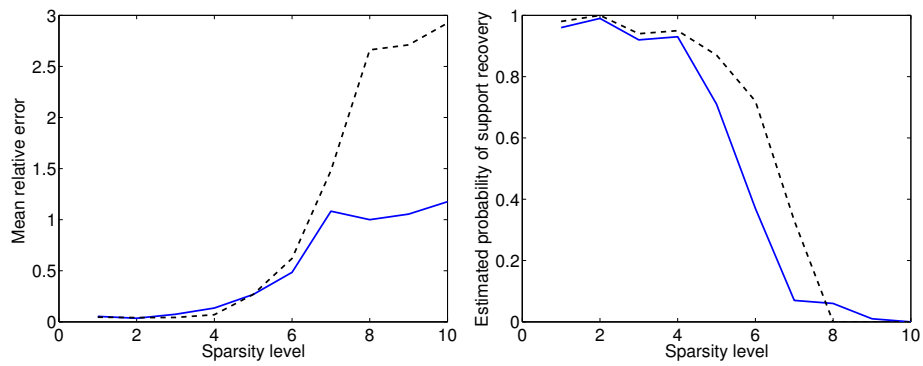


Figure 2: Mean relative error (left) and estimated probability of support recovery (right) in the noisy case for the convex relaxation (13) (plain line) and the greedy method applied to (12) (dashed line).

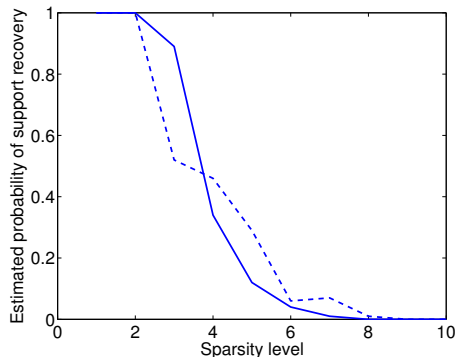


Figure 3: Estimated probability of support recovery from the magnitude of the Fourier transform when using (35) for  $N = n$  (plain line) and when using (37) for  $N = 2n$  (dashed line) versus the sparsity level.

**Support estimation from the power spectrum.** We now test if the convex method is robust to invariances, such as the ones discussed in Sect. 5.2, when estimating the support. In particular, we start with a setting in which  $\mathbf{x} \in \mathbb{R}^n$ ,  $N = n = 20$  and  $\mathbf{q}_i^H$ ,  $i = 1, \dots, N$ , are the rows of the  $n$ -point Fourier matrix. In this case, the autocorrelation cannot be computed and the convex formulation (35) is used to estimate the support. Then, we perform similar experiments but with oversampling ( $N = 2n$ ), thus allowing for the computation of the autocorrelation and the use of (37). Results shown in Fig. 3 indicate that by using (35) we can recover the correct support for sufficiently sparse signals without oversampling, i.e., with as many measurements as unknowns, while using more information extracted from the autocorrelation of the signal only slightly helps to recover larger supports.

## 7 Conclusions

The paper proposed a new approach to phase retrieval of sparse signals. This approach is based on a group-sparse optimization formulation of the problem with linearized constraints. Exact and stable recovery results were shown for a convex relaxation of this formulation both in the real and complex case. Invariances to circular shifts and reflections that are common in phase retrieval problems were also discussed and a practical technique was given to prevent these from breaking the sparsity of the solution.

Future work will focus on deriving theoretical guarantees for the invariance issue. Another direction of research concerns the analysis of greedy algorithms for the group-sparse optimization problem, which proved as valuable as the convex relaxations in experiments for measurements without invariance. How to deal with invariances in these methods will also be investigated.

## References

- [1] E. J. Candès. Compressive sampling. In *Proceedings of the International Congress of Mathematicians: invited lectures*, pages 1433–1452, 2006.
- [2] E. J. Candès, Y. C. Eldar, T. Strohmer, and V. Voroninski. Phase retrieval via matrix completion. *SIAM Journal on Imaging Sciences*, 6(1):199–225, 2013.
- [3] E. J. Candès, T. Strohmer, and V. Voroninski. Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming. *Communications on Pure and Applied Mathematics*, 66(8):1241–1274, 2013.
- [4] E. J. Candès, M. B. Wakin, and S. P. Boyd. Enhancing sparsity by reweighted  $\ell_1$  minimization. *Journal of Fourier Analysis and Applications*, 14(5):877–905, 2008.

- [5] D. L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):1289–1306, 2006.
- [6] D. L. Donoho, M. Elad, and V. N. Temlyakov. Stable recovery of sparse overcomplete representations in the presence of noise. *IEEE Transactions on Information Theory*, 52(1):6–18, 2006.
- [7] M. Ehler, M. Fornasier, and J. Sigl. Quasi-linear compressed sensing. Technical report, 2013. [http://www-m15.ma.tum.de/foswiki/pub/M15/Allgemeines/PublicationsEN/greedy\\_21.pdf](http://www-m15.ma.tum.de/foswiki/pub/M15/Allgemeines/PublicationsEN/greedy_21.pdf).
- [8] J. Fienup. Phase retrieval algorithms: a comparison. *Applied Optics*, 21(15):2758–2769, 1982.
- [9] R. Gerchberg and W. Saxton. A practical algorithm for the determination of phase from image and diffraction plane pictures. *Optik*, 35:237–246, 1972.
- [10] R. Gonsalves. Phase retrieval from modulus data. *Journal of Optical Society of America*, 66(9):961–964, 1976.
- [11] R. W. Harrison. Phase problem in crystallography. *Journal of the Optical Society of America A*, 10(5):1046–1055, 1993.
- [12] K. Jaganathan, S. Oymak, and B. Hassibi. Recovery of sparse 1-d signals from the magnitudes of their Fourier transform. In *IEEE International Symposium on Information Theory (ISIT)*, pages 1473–1477, 2012.
- [13] D. Kohler and L. Mandel. Source reconstruction from the modulus of the correlation function: a practical approach to the phase problem of optical coherence theory. *Journal of the Optical Society of America*, 63(2):126–134, 1973.
- [14] F. Lauer and H. Ohlsson. Finding sparse solutions of systems of polynomial equations via group-sparsity optimization. *arXiv preprint*, arXiv:1311.5871, 2013.
- [15] S. Mukherjee and C. S. Seelamantula. An iterative algorithm for phase retrieval with sparsity constraints: application to frequency domain optical coherence tomography. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 553–556, 2012.
- [16] H. Ohlsson and Y. C. Eldar. On conditions for uniqueness in sparse phase retrieval. *CoRR*, abs/1308.5447, 2013.
- [17] H. Ohlsson, A. Y. Yang, R. Dong, and S. Sastry. Nonlinear basis pursuit. *arXiv preprint arXiv:1304.5802*, 2013.
- [18] H. Ohlsson, A. Y. Yang, R. Dong, M. Verhaegen, and S. Sastry. Quadratic basis pursuit. *arXiv preprint arXiv:1301.7002*, 2013.
- [19] J. Ranieri, A. Chebira, Y. M. Lu, and M. Vetterli. Phase retrieval for sparse signals: Uniqueness conditions. *CoRR*, abs/1308.3058, 2013.
- [20] Y. Shechtman, A. Beck, and Y. C. Eldar. GESPAR: Efficient phase retrieval of sparse signals. *arXiv preprint arXiv:1301.1018*, 2013.
- [21] Y. Shechtman, Y. C. Eldar, A. Szameit, and M. Segev. Sparsity based sub-wavelength imaging with partially incoherent light via quadratic compressed sensing. *Optics Express*, 19(16):14807–14822, 2011.
- [22] I. Waldspurger, A. d’Aspremont, and S. Mallat. Phase recovery, maxcut and complex semidefinite programming. *arXiv preprint arXiv:1206.0102*, 2012.
- [23] A. Walther. The question of phase retrieval in optics. *Journal of Modern Optics*, 10(1):41–49, 1963.

## A Lemmas

**Lemma 1.** Let  $\nu$  and  $\nu^{-1}$  be defined as in Section 2. For all  $\mathbf{v} \in \mathbb{R}^M$  such that  $\nu^{-1}(\mathbf{v}) \neq \mathbf{0}$  and  $\mathbf{v}' \in \mathbb{R}^M$ , if  $(\nu(\nu^{-1}(\mathbf{v})))_{jj} \neq v'_{jj}$  for some  $j \in \{1, \dots, n\}$ , then  $\mathbf{v} \neq \mathbf{v}'$ .

*Proof.* Assume that  $\nu^{-1}(\mathbf{v}) \neq \mathbf{0}$  and let  $i = \min_j j$ , s.t.  $v_{jj} > 0$ . Then,  $\forall j \in \{1, \dots, n\}$ , we have

$$(\nu(\nu^{-1}(\mathbf{v})))_{jj} = (\nu^{-1}(\mathbf{v}))_j^2 = \frac{v_{ji}^2}{v_{ii}} = v_{jj},$$

where the last equality follows from the definition of  $\nu^{-1}$ . Therefore, if  $(\nu(\nu^{-1}(\mathbf{v})))_{jj} = v_{jj} \neq v'_{jj}$  for some  $j \in \{1, \dots, n\}$ , we obtain  $\mathbf{v} \neq \mathbf{v}'$ .  $\square$

**Lemma 2.** Let  $\mathbf{A} = [\mathbf{A}_1, \dots, \mathbf{A}_M]$  be an  $N \times M$  real matrix with mutual coherence  $\mu(\mathbf{A})$  as in Definition 1. Let  $\mathbf{W}$  be the  $M \times M$ -diagonal matrix of entries  $w_i = \|\mathbf{A}_i\|_2$ . Then, for all  $\boldsymbol{\delta} \in \text{Ker}(\mathbf{A})$  and  $i \in \{1, \dots, M\}$ , the bound

$$w_i^2 \delta_i^2 \leq \frac{\mu^2(\mathbf{A})}{1 + \mu^2(\mathbf{A})} \|\mathbf{W}\boldsymbol{\delta}\|_2^2 \quad (38)$$

holds.

*Proof.* See Lemma 2 in [14].  $\square$

**Lemma 3.** Let the complex Veronese map  $\nu$  be as in Definition 2. Then, for all  $\mathbf{q} \in \mathbb{C}^n$  and  $\mathbf{x} \in \mathbb{C}^n$ , the following equality holds:

$$\mathbf{x}^H \mathbf{q} \mathbf{q}^H \mathbf{x} = 2\Re(\nu(\mathbf{q})^H \nu(\mathbf{x})) - \sum_{j=1}^n (\nu(\mathbf{q}))_{jj} (\nu(\mathbf{x}))_{jj}.$$

*Proof.*

$$\begin{aligned} \mathbf{x}^H \mathbf{q} \mathbf{q}^H \mathbf{x} &= \left( \sum_{k=1}^n \bar{x}_k q_k \right) \left( \sum_{j=1}^n \bar{q}_j x_j \right) \\ &= \sum_{j=1}^n \bar{q}_j x_j \sum_{k=1}^n q_k \bar{x}_k = \sum_{j=1}^n \sum_{k=1}^n \bar{q}_j q_k x_j \bar{x}_k \\ &= \sum_{j=1}^n q_j \bar{q}_j x_j \bar{x}_j + \sum_{j=1}^n \sum_{k \neq j} \bar{q}_j q_k x_j \bar{x}_k \\ &= \sum_{j=1}^n q_j \bar{q}_j x_j \bar{x}_j + \sum_{j=1}^n \left( \sum_{k < j} \bar{q}_j q_k x_j \bar{x}_k + \sum_{k > j} \bar{q}_j q_k x_j \bar{x}_k \right) \\ &= \sum_{j=1}^n q_j \bar{q}_j x_j \bar{x}_j + \sum_{j=1}^n \sum_{k > j} (\bar{q}_k q_j x_k \bar{x}_j + \bar{q}_j q_k x_j \bar{x}_k) \\ &= \sum_{j=1}^n q_j \bar{q}_j x_j \bar{x}_j + \sum_{j=1}^n \sum_{k > j} (\overline{\bar{q}_j q_k x_j \bar{x}_k} + \bar{q}_j q_k x_j \bar{x}_k) \end{aligned}$$

At this point, we use the fact that  $z + \bar{z} = 2\Re(z)$ , which yields

$$\begin{aligned} \mathbf{x}^H \mathbf{q} \mathbf{q}^H \mathbf{x} &= \sum_{j=1}^n q_j \bar{q}_j x_j \bar{x}_j + 2 \sum_{j=1}^n \sum_{k > j} \Re(\bar{q}_j q_k x_j \bar{x}_k) \\ &= \sum_{j=1}^n (\nu(\mathbf{q}))_{jj} (\nu(\mathbf{x}))_{jj} + 2 \sum_{j=1}^n \sum_{k > j} \Re\left( (\nu(\mathbf{q}))_{jk} (\nu(\mathbf{x}))_{jk} \right) \end{aligned}$$

Since  $(\nu(\mathbf{q}))_{jj}(\nu(\mathbf{x}))_{jj} = q_j \bar{q}_j x_j \bar{x}_j = |q_j x_j|^2$  is a real number, it can be introduced in the second sum as

$$\begin{aligned}
\mathbf{x}^H \mathbf{q} \mathbf{q}^H \mathbf{x} &= - \sum_{j=1}^n (\nu(\mathbf{q}))_{jj} (\nu(\mathbf{x}))_{jj} + 2 \sum_{j=1}^n \sum_{k \geq j} \Re \left( \left( \overline{\nu(\mathbf{q})} \right)_{jk} (\nu(\mathbf{x}))_{jk} \right) \\
&= - \sum_{j=1}^n (\nu(\mathbf{q}))_{jj} (\nu(\mathbf{x}))_{jj} + 2 \sum_{j=1}^n \sum_{k \geq j} \left[ \Re \left( \overline{\nu(\mathbf{q})} \right)_{jk} \Re(\nu(\mathbf{x}))_{jk} - \Im \left( \overline{\nu(\mathbf{q})} \right)_{jk} \Im(\nu(\mathbf{x}))_{jk} \right] \\
&= - \sum_{j=1}^n (\nu(\mathbf{q}))_{jj} (\nu(\mathbf{x}))_{jj} + 2 \left[ \Re \left( \overline{\nu(\mathbf{q})} \right)^T \Re(\nu(\mathbf{x})) - \Im \left( \overline{\nu(\mathbf{q})} \right)^T \Im(\nu(\mathbf{x})) \right] \\
&= - \sum_{j=1}^n (\nu(\mathbf{q}))_{jj} (\nu(\mathbf{x}))_{jj} + 2 \left[ \Re(\nu(\mathbf{q}))^T \Re(\nu(\mathbf{x})) + \Im(\nu(\mathbf{q}))^T \Im(\nu(\mathbf{x})) \right] \\
&= 2 \Re(\nu(\mathbf{q})^H \nu(\mathbf{x})) - \sum_{j=1}^n (\nu(\mathbf{q}))_{jj} (\nu(\mathbf{x}))_{jj},
\end{aligned}$$

where the last equality is due to  $\Re(\mathbf{a}^H \mathbf{b}) = \Re(\mathbf{a})^T \Re(\mathbf{b}) + \Im(\mathbf{a})^T \Im(\mathbf{b})$ .  $\square$

**Lemma 4.** Let  $\nu^{-1}$  be as in Definition 3 and the matrices  $\mathbf{W}_j$  as in Sect. 3.1. Then, for all  $\mathbf{v} \in \mathbb{C}^M$ ,

$$\|\nu^{-1}(\mathbf{v})\|_0 \leq \|\{\mathbf{W}_j \mathbf{v}\}_{j=1}^n\|_0. \quad (39)$$

*Proof.* According to Definition 3,  $|(\nu^{-1}(\mathbf{v}))_j|^2 = \frac{|v_{ji}|^2}{v_{ii}} = v_{jj}$ . On the other hand,  $v_{jj}$  belongs to a single group of variables generated by the  $\mathbf{W}_j$ , i.e.,  $(\mathbf{W}_k)_{(jj)} \neq \mathbf{0} \Leftrightarrow k = j$ . Thus,  $(\nu^{-1}(\mathbf{v}))_j \neq 0 \Rightarrow v_{jj} > 0 \Rightarrow \mathbf{W}_j \mathbf{v} \neq \mathbf{0} \Rightarrow \|\nu^{-1}(\mathbf{v})\|_0 \leq \|\{\mathbf{W}_j \mathbf{v}\}_{j=1}^n\|_0$ .  $\square$

However, note that the converse is not true: we can have  $\|\nu^{-1}(\mathbf{v})\|_0 \neq \|\{\mathbf{W}_j \mathbf{v}\}_{j=1}^n\|_0$ , e.g., when  $v_{jj} = 0$ ,  $j = 1, \dots, n$ , and  $v_{ij} \neq 0$  for some  $i$  and  $j$ .

**Lemma 5.** Let  $\nu$  and  $\nu^{-1}$  be defined as in Definitions 2 and 3. For all  $\mathbf{v} \in \mathbb{C}^M$  such that  $\nu^{-1}(\mathbf{v}) \neq \mathbf{0}$  and  $\mathbf{v}' \in \mathbb{C}^M$ , if  $(\nu(\nu^{-1}(\mathbf{v})))_{jj} \neq v'_{jj}$  for some  $j \in \{1, \dots, n\}$ , then  $\mathbf{v} \neq \mathbf{v}'$ .

*Proof.* Assume that  $\nu^{-1}(\mathbf{v}) \neq \mathbf{0}$  and let  $i = \min_j j$ , s.t.  $\Re(v_{jj}) > 0$ ,  $\Im(v_{jj}) = 0$ . Then,  $\forall j \in \{1, \dots, n\}$ , we have

$$(\nu(\nu^{-1}(\mathbf{v})))_{jj} = (\nu^{-1}(\mathbf{v}))_j \left( \overline{\nu^{-1}(\mathbf{v})} \right)_j = \frac{\bar{v}_{ij}}{\sqrt{v_{ii}}} \frac{v_{ij}}{\sqrt{v_{ii}}} = \frac{|v_{ij}|^2}{v_{ii}} = v_{jj},$$

where the last equality follows from Definition 3. Therefore, if  $(\nu(\nu^{-1}(\mathbf{v})))_{jj} = v_{jj} \neq v'_{jj}$  for some  $j \in \{1, \dots, n\}$ , we obtain  $\mathbf{v} \neq \mathbf{v}'$ .  $\square$

**Lemma 6.** For all  $\delta \in \mathbb{C}^M$  such that  $\Re(\mathbf{A}\delta) = \mathbf{0}$ , the inequality

$$(w_i^R)^2 \Re(\delta_i)^2 + (w_i^I)^2 \Im(\delta_i)^2 \leq \frac{2\mu^2(\tilde{\mathbf{A}})}{1 + \mu^2(\tilde{\mathbf{A}})} \|\mathbf{W}^R \Re(\delta) + i \mathbf{W}^I \Im(\delta)\|^2$$

holds with  $\tilde{\mathbf{A}} = [\Re(\mathbf{A}), -\Im(\mathbf{A})] \in \mathbb{R}^{N \times 2M}$ .

*Proof.* We rewrite the assumption as

$$\mathbf{0} = \Re(\mathbf{A}\delta) = \tilde{\mathbf{A}}\tilde{\delta},$$

where  $\tilde{\delta} = [\Re(\delta)^T, \Im(\delta)^T]^T \in \mathbb{R}^{2M}$ . Then, we can use Lemma 2 to bound the entries in  $\tilde{\delta}$  as

$$(w_i^R)^2 \Re(\delta_i)^2 \leq \frac{\mu^2(\tilde{\mathbf{A}})}{1 + \mu^2(\tilde{\mathbf{A}})} \|\tilde{\mathbf{W}}\tilde{\delta}\|_2^2,$$

with  $\tilde{\mathbf{W}} = \begin{pmatrix} \mathbf{W}^R & \mathbf{0} \\ \mathbf{0} & \mathbf{W}^I \end{pmatrix}$  a diagonal matrix of entries  $(\tilde{\mathbf{W}})_{i,i} = \|\tilde{\mathbf{A}}_i\|_2$ , and

$$(w_i^I)^2 \Im(\delta_i)^2 \leq \frac{\mu^2(\tilde{\mathbf{A}})}{1 + \mu^2(\tilde{\mathbf{A}})} \|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2.$$

These inequalities lead to

$$(w_i^R)^2 \Re(\delta_i)^2 + (w_i^I)^2 \Im(\delta_i)^2 \leq \frac{2\mu^2(\tilde{\mathbf{A}})}{1 + \mu^2(\tilde{\mathbf{A}})} \|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2,$$

where

$$\|\tilde{\mathbf{W}}\tilde{\boldsymbol{\delta}}\|_2^2 = \|\mathbf{W}^R \Re(\boldsymbol{\delta})\|_2^2 + \|\mathbf{W}^I \Im(\boldsymbol{\delta})\|_2^2 = \|\mathbf{W}^R \Re(\boldsymbol{\delta}) + i\mathbf{W}^I \Im(\boldsymbol{\delta})\|_2^2.$$

□

**Lemma 7.** For all  $\boldsymbol{\delta} \in \mathbb{R}^M$  such that  $\Re(\mathbf{A})\boldsymbol{\delta} = \mathbf{0}$ , the inequality

$$(w_i^R)^2 \delta_i^2 \leq \frac{\mu^2(\Re(\mathbf{A}))}{1 + \mu^2(\Re(\mathbf{A}))} \|\mathbf{W}^R \boldsymbol{\delta}\|_2^2,$$

where  $\mathbf{W}^R$  is a diagonal matrix of entries  $(\mathbf{W}^R)_{i,i} = w_i^R = \|\Re(\mathbf{A}_i)\|_2$ .

*Proof.* This is a direct consequence of Lemma 2 applied to  $\Re(\mathbf{A})$  and  $\boldsymbol{\delta} \in \text{Ker}(\Re(\mathbf{A}))$ . □

## B Proofs

### B.1 Proof of Theorem 2

This proof is similar to the one of Theorem 3 in [14].

*Proof.* The vector  $\mathbf{v}_0$  is the unique solution to (6) if the inequality

$$\sum_{j=1}^n \|\mathbf{W}_j \mathbf{W}(\mathbf{v}_0 + \boldsymbol{\delta})\|_2 > \sum_{j=1}^n \|\mathbf{W}_j \mathbf{W} \mathbf{v}_0\|_2$$

holds for all  $\boldsymbol{\delta} \neq \mathbf{0}$  satisfying  $\mathbf{A}\boldsymbol{\delta} = \mathbf{0}$ . The inequality above can be rewritten as

$$\sum_{j \in I_0} \|\mathbf{W}_j \mathbf{W} \boldsymbol{\delta}\|_2 + \sum_{j \notin I_0} \|\mathbf{W}_j \mathbf{W}(\mathbf{v}_0 + \boldsymbol{\delta})\|_2 - \|\mathbf{W}_j \mathbf{W} \mathbf{v}_0\|_2 > 0,$$

where  $I_0 = \{j \in \{1, \dots, n\} : \mathbf{W}_j \mathbf{W} \mathbf{v}_0 = \mathbf{0}\}$ . By the triangle inequality,  $\|\mathbf{a} + \mathbf{b}\|_2 - \|\mathbf{a}\|_2 \geq -\|\mathbf{b}\|_2$ , this condition is met if

$$\sum_{j \in I_0} \|\mathbf{W}_j \mathbf{W} \boldsymbol{\delta}\|_2 - \sum_{j \notin I_0} \|\mathbf{W}_j \mathbf{W} \boldsymbol{\delta}\|_2 > 0$$

or

$$\sum_{j=1}^n \|\mathbf{W}_j \mathbf{W} \boldsymbol{\delta}\|_2 - 2 \sum_{j \notin I_0} \|\mathbf{W}_j \mathbf{W} \boldsymbol{\delta}\|_2 > 0. \quad (40)$$

By defining  $G_j$  as the set of indexes corresponding to nonzero columns of  $\mathbf{W}_j$ , Lemma 2 yields

$$\|\mathbf{W}_j \mathbf{W} \boldsymbol{\delta}\|_2^2 = \sum_{i \in G_j} w_i^2 \delta_i^2 \leq n \frac{\mu^2(\mathbf{A})}{1 + \mu^2(\mathbf{A})} \|\mathbf{W} \boldsymbol{\delta}\|_2^2,$$

Due to the fact that  $\bigcup_{k \in \{1, \dots, n\}} G_k = \{1, \dots, M\}$ , we also have

$$\|\mathbf{W} \boldsymbol{\delta}\|_2^2 = \sum_{i=1}^M w_i^2 \delta_i^2 \leq \sum_{k=1}^n \sum_{i \in G_k} w_i^2 \delta_i^2 = \sum_{k=1}^n \|\mathbf{W}_k \mathbf{W} \boldsymbol{\delta}\|_2^2 \leq \left( \sum_{k=1}^n \|\mathbf{W}_k \mathbf{W} \boldsymbol{\delta}\|_2 \right)^2,$$

which then leads to

$$\|\mathbf{W}_j \mathbf{W} \boldsymbol{\delta}\|_2^2 \leq n \frac{\mu^2(\mathbf{A})}{1 + \mu^2(\mathbf{A})} \left( \sum_{k=1}^n \|\mathbf{W}_k \mathbf{W} \boldsymbol{\delta}\|_2 \right)^2.$$

Introducing this result in (40) gives the condition

$$\sum_{j=1}^n \|\mathbf{W}_j \mathbf{W} \boldsymbol{\delta}\|_2 - 2(n - |I_0|) \frac{\mu(\mathbf{A})\sqrt{n}}{\sqrt{1 + \mu^2(\mathbf{A})}} \sum_{k=1}^n \|\mathbf{W}_k \mathbf{W} \boldsymbol{\delta}\|_2 > 0.$$

Finally, given that  $|I_0| = n - \|\{\mathbf{W}_j \mathbf{v}_0\}_{j=1}^n\|_0 = n - \|\mathbf{x}_0\|_0$ , this yields

$$\sum_{j=1}^n \|\mathbf{W}_j \mathbf{W} \boldsymbol{\delta}\|_2 - 2\|\mathbf{x}_0\|_0 \frac{\mu(\mathbf{A})\sqrt{n}}{\sqrt{1 + \mu^2(\mathbf{A})}} \sum_{k=1}^n \|\mathbf{W}_k \mathbf{W} \boldsymbol{\delta}\|_2 > 0.$$

or, after rearranging the terms,

$$\|\mathbf{x}_0\|_0 < \frac{\sqrt{1 + \mu^2(\mathbf{A})}}{2\mu(\mathbf{A})\sqrt{n}},$$

which can be rewritten as in the statement of the Theorem.  $\square$

## B.2 Proof of Theorem 7

This proof is very similar to the ones of Theorems 2 and 4.

*Proof.* The vector  $\mathbf{v}_0$  is the unique solution to (31) if the inequality

$$\sum_{j=1}^n \|\mathbf{W}_j^R(\mathbf{v}_0 + \boldsymbol{\delta})\|_2 > \sum_{j=1}^n \|\mathbf{W}_j^R \mathbf{v}_0\|_2$$

holds for all  $\boldsymbol{\delta} \in \mathbb{R}^M$  such that  $\Re(\mathbf{A})(\mathbf{v}_0 + \boldsymbol{\delta}) = \mathbf{y}$ , which implies the constraint  $\Re(\mathbf{A})\boldsymbol{\delta} = \mathbf{0}$  on  $\boldsymbol{\delta}$ . The inequality above can be rewritten as

$$\sum_{j \in I_0} \|\mathbf{W}_j^R \boldsymbol{\delta}\|_2 + \sum_{j \notin I_0} \|\mathbf{W}_j^R(\mathbf{v}_0 + \boldsymbol{\delta})\|_2 - \sum_{j \notin I_0} \|\mathbf{W}_j^R \mathbf{v}_0\|_2 > 0,$$

where  $I_0 = \{j \in \{1, \dots, n\} : \mathbf{W}_j^R \mathbf{v}_0 = \mathbf{0}\}$ . By the triangle inequality,  $\|\mathbf{a} + \mathbf{b}\| - \|\mathbf{a}\| \geq -\|\mathbf{b}\|$  with  $\mathbf{a} = \mathbf{W}_j^R \mathbf{v}_0$ , this condition is met if

$$\sum_{j=1}^n \|\mathbf{W}_j^R \boldsymbol{\delta}\|_2 - 2 \sum_{j \notin I_0} \|\mathbf{W}_j^R \boldsymbol{\delta}\|_2 > 0. \quad (41)$$

By defining  $G_j$  as the set of indexes corresponding to nonzero columns of  $\mathbf{W}_j$ , Lemma 7 yields

$$\|\mathbf{W}_j^R \boldsymbol{\delta}\|_2^2 = \sum_{i \in G_j} (w_i^R)^2 \delta_i^2 \leq n \frac{\mu^2(\Re(\mathbf{A}))}{1 + \mu^2(\Re(\mathbf{A}))} \|\mathbf{W}^R \boldsymbol{\delta}\|_2^2.$$

Due to the fact that  $\bigcup_{k \in \{1, \dots, n\}} G_k = \{1, \dots, M\}$ , we also have

$$\|\mathbf{W}^R \boldsymbol{\delta}\|_2^2 = \sum_{i=1}^M (w_i^R)^2 \delta_i^2 \leq \sum_{k=1}^n \sum_{i \in G_k} (w_i^R)^2 \delta_i^2 = \sum_{k=1}^n \|\mathbf{W}_k^R \boldsymbol{\delta}\|_2^2 \leq \left( \sum_{k=1}^n \|\mathbf{W}_k^R \boldsymbol{\delta}\|_2 \right)^2,$$

which then leads to

$$\|\mathbf{W}_j^R \boldsymbol{\delta}\|_2^2 \leq n \frac{\mu^2(\Re(\mathbf{A}))}{1 + \mu^2(\Re(\mathbf{A}))} \left( \sum_{k=1}^n \|\mathbf{W}_k^R \boldsymbol{\delta}\|_2 \right)^2.$$

Introducing this result in (41) gives the condition

$$\sum_{j=1}^n \|\mathbf{W}_j^R \boldsymbol{\delta}\|_2 - 2(n - |I_0|) \frac{\mu(\Re(\mathbf{A}))\sqrt{n}}{\sqrt{1 + \mu^2(\Re(\mathbf{A}))}} \sum_{k=1}^n \|\mathbf{W}_k^R \boldsymbol{\delta}\|_2 > 0.$$

Finally, given that  $|I_0| = n - \|\mathbf{x}_0\|_0$ , this yields, for  $\boldsymbol{\delta} \neq \mathbf{0}$ ,

$$\|\mathbf{x}_0\|_0 < \frac{\sqrt{1 + \mu^2(\Re(\mathbf{A}))}}{2\mu(\Re(\mathbf{A}))\sqrt{n}},$$

which can be rewritten as in the statement of the Theorem.  $\square$

### B.3 Proof of Proposition 1

*Proof.* To prove Statement 1, note that the operation  $\mathbf{x}_2 = \text{shift}(\text{reflection}(\mathbf{x}_1), 1)$  is equivalent to

$$\mathbf{x}_2 = \begin{bmatrix} x_{11} \\ \text{reflection}(\tilde{\mathbf{x}}_1) \end{bmatrix}, \quad \tilde{\mathbf{x}}_1 = [x_{12}, x_{13}, \dots, x_{1n}]^T,$$

where the reflection is centered on  $\tilde{x}_{1(n/2)} = x_{1(1+n/2)}$ . Thus,

$$\mathbf{x}_2 = [x_{11}, x_{1n}, x_{1(n-1)}, \dots, x_{12}]^T$$

and

$$\sum_{i=2}^{n/2} |x_{2i}|^2 = \sum_{i=0}^{n/2-2} |x_{1(n-i)}|^2 = \sum_{i=2+\frac{n}{2}}^n |x_{1i}|^2 > \sum_{i=2}^{n/2} |x_{1i}|^2 = \sum_{i=0}^{n/2-2} |x_{2(n-i)}|^2 = \sum_{i=2+\frac{n}{2}}^n |x_{2i}|^2,$$

where the inequality holds whenever this operation is performed in (34) (i.e., if this is not the case, then  $\mathbf{x}_2 = \mathbf{x}_1$ ). Therefore, for all  $\mathbf{x}$ , we obtain a vector  $\mathbf{x}_2 = \varphi(\mathbf{x})$  such that  $\sum_{i=2}^{n/2} |x_{2i}|^2 \geq \sum_{i=2+\frac{n}{2}}^n |x_{2i}|^2$ . Since  $\arg \max_{i \in \{1, \dots, n\}} |x_{2i}| = 1$  and  $\text{shift}(\mathbf{x}_2, 1 - \arg \max_{i \in \{1, \dots, n\}} |x_{2i}|) = \mathbf{x}_2$ , this implies that  $\varphi(\mathbf{x}_2) = \mathbf{x}_2$ , i.e.,  $\varphi(\varphi(\mathbf{x})) = \varphi(\mathbf{x})$ .

Statement 2 is easily seen from the fact that  $\mathbf{x}_1$  does not change when  $\mathbf{x}$  is shifted.

To prove Statement 3, let  $\mathbf{x}' = \text{reflection}(\mathbf{x}) = [x_n, x_{n-1}, \dots, x_1]^T$ ,  $k = \arg \max_{i \in \{1, \dots, n\}} |x_i|$ , and  $k' = \arg \max_{i \in \{1, \dots, n\}} |x'_i|$ . Then,  $k' = n - k + 1$  and

$$\begin{aligned} \mathbf{x}'_1 &= \text{shift}(\mathbf{x}', 1 - k') = \text{shift}(\mathbf{x}', k - n) = \text{shift}(\mathbf{x}', k) \\ &= [x'_{n-k+1}, x'_{n-k+2}, \dots, x'_n, x'_1, x'_2, \dots, x'_{n-k}]^T \\ &= [x_k, x_{k-1}, \dots, x_1, x_n, x_{n-1}, \dots, x_{k+1}]^T. \end{aligned}$$

Define  $\mathbf{x}_1 = \text{shift}(\mathbf{x}, 1 - k) = [x_k, x_{k+1}, \dots, x_n, x_1, x_2, \dots, x_{k-1}]^T$ .

If  $k < n/2$ , we have

$$\sum_{i=2}^{n/2} |x'_{1i}|^2 = \sum_{i=1}^{k-1} |x_i|^2 + \sum_{i=1}^{\frac{n}{2}-k} |x_{n-i+1}|^2 = \sum_{i=n/2-k}^n |x_{1i}|^2 + \sum_{i=n/2+2}^{n-k+1} |x_{1i}|^2 = \sum_{i=n/2+2}^n |x_{1i}|^2,$$

and otherwise, if  $k \geq n/2$ , we obtain

$$\sum_{i=2}^{n/2} |x'_{1i}|^2 = \sum_{i=k-\frac{n}{2}+1}^{k-1} |x_i|^2 = \sum_{i=n/2+2}^n |x_{1i}|^2.$$

On the other hand, if  $k \leq n/2$ , then

$$\sum_{i=2+\frac{n}{2}}^n |x'_{1i}|^2 = \sum_{i=k+1}^{\frac{n}{2}+k-1} |x_i|^2 = \sum_{i=2}^{n/2} |x_{1i}|^2,$$



while if  $k > n/2$ :

$$\sum_{i=2+\frac{n}{2}}^n |x'_{1i}|^2 = \sum_{i=k+1}^n |x_i|^2 + \sum_{i=1}^{k-\frac{n}{2}-1} |x_i|^2 = \sum_{i=2}^{n-k+1} |x_{1i}|^2 + \sum_{i=n-k+2}^{n/2} |x_{1i}|^2 = \sum_{i=2}^{n/2} |x_{1i}|^2.$$

Therefore,

$$\sum_{i=2}^{n/2} |x_{1i}|^2 > \sum_{i=2+\frac{n}{2}}^n |x_{1i}|^2 \Leftrightarrow \sum_{i=2}^{n/2} |x'_{1i}|^2 < \sum_{i=2+\frac{n}{2}}^n |x'_{1i}|^2$$

and  $\varphi$  applies a reflection to  $\mathbf{x}$  if and only if it does not apply one to  $\mathbf{x}'$  (the inequalities above are strict by assumption). Thus, in the case  $\mathbf{x}$  is reflected by  $\varphi$ , we have

$$\varphi(\mathbf{x}') = \text{shift}(\text{reflection}(\mathbf{x}), k) = \text{reflection}(\text{shift}(\mathbf{x}, n - k))$$

and

$$\begin{aligned} \varphi(\mathbf{x}) &= \text{shift}(\text{reflection}(\text{shift}(\mathbf{x}, 1 - k)), 1) \\ &= \text{reflection}(\text{shift}(\text{shift}(\mathbf{x}, 1 - k), n - 1)) \\ &= \text{reflection}(\text{shift}(\mathbf{x}, n - k)) \\ &= \varphi(\mathbf{x}'), \end{aligned}$$

while in the case  $\mathbf{x}$  is not reflected by  $\varphi$ , we have

$$\varphi(\mathbf{x}) = \text{shift}(\mathbf{x}, 1 - k)$$

and

$$\begin{aligned} \varphi(\mathbf{x}') &= \text{shift}(\text{reflection}(\text{shift}(\text{reflection}(\mathbf{x}), k)), 1) \\ &= \text{reflection}(\text{shift}(\text{shift}(\text{reflection}(\mathbf{x}), k), n - 1)) \\ &= \text{reflection}(\text{shift}(\text{reflection}(\mathbf{x}), n + k - 1)) \\ &= \text{reflection}(\text{reflection}(\text{shift}(\mathbf{x}, -k + 1))) \\ &= \text{shift}(\mathbf{x}, 1 - k) \\ &= \varphi(\mathbf{x}). \end{aligned}$$

□