

# Role of Sparsity and Structure in the Optimization Landscape of Non-convex Matrix Sensing

Igor Molybog · Somayeh Sojoudi · Javad Lavaei

Received: date / Accepted: date

**Abstract** In this work, we study the optimization landscape of the non-convex matrix sensing problem that is known to have many local minima in the worst case. Since the existing results are related to the notion of restricted isometry property (RIP) that cannot directly capture the underlying structure of a given problem, they can hardly be applied to real-world problems where the amount of data is not exorbitantly high. To address this issue, we develop the notion of kernel structure property to obtain necessary and sufficient conditions for the inexistence of spurious local solution of any class of matrix sensing problems over a given search space. This notion precisely captures the underlying sparsity and structure of the problem, based on tools in conic optimization. We simplify the conditions for a certain class of problems to show their satisfaction and apply them to data analytics for power systems.

**Keywords** Non-convex optimization · spurious local minima · matrix sensing

## 1 Introduction

Even under the ideal condition of no noise and zero approximation error, many highly-efficient machine learning techniques involve solving potentially hard or intractable computational problems while learning from data. In practice, they are tackled by heuristic optimization algorithms, based on relaxations or greedy principals.

---

Igor Molybog  
University of California, Berkeley, CA  
E-mail: igormolybog@berkeley.edu

Somayeh Sojoudi  
University of California, Berkeley, CA  
E-mail: sojoudi@berkeley.edu

Javad Lavaei (corresponding author)  
University of California, Berkeley, CA  
E-mail: lavaei@berkeley.edu

The lack of guarantees on their performance limits their use in applications with significant cost of an error, impacting our ability to implement progressive data analysis techniques in crucial social and economic systems, such as healthcare, transportation, and energy production and distribution. Commonly, non-convexity is the main obstacle for a guaranteed learning of continuous parameters.

It is well known that many fundamental problems with a natural non-convex formulation can be  $\mathcal{NP}$ -hard (Pardalos and Vavasis 1991). Sophisticated techniques for addressing this issue, like generic convex relaxations, may require working in an unrealistically high-dimensional space to guarantee exactness of the solution. As a consequence of complicated geometrical structures, a non-convex function may contain an exponential number of saddle points and spurious local minima, and therefore local search algorithms may become trapped in such points. Nevertheless, empirical observations show positive results regarding the application of these approaches to several practically important instances. This provokes a large branch of research that aims to explain the success of experimental results in order to understand the boundaries of applicability of the existing algorithms and develop new ones. A recent direction in non-convex optimization consists in studying how simple algorithms can solve potentially hard problems arising in data analysis applications. The most commonly applied class of such algorithms is based on *local search*, which will be the focus of this work. In some cases, prior information about the location of the solution is available, which significantly reduces the complexity of the search.

Consider searching over some given domain  $\mathcal{X}$ . For a twice continuously differentiable objective function  $f: \mathcal{X} \rightarrow \mathbb{R}$  that reaches its global minimum  $f^*$ , if the point  $x$  attains  $f(x) = f^*$ , then we call it a *global minimizer*. The point  $x$  is said to be a *local minimizer* if  $f(x) \leq f(x')$  holds for all  $x'$  within a local neighborhood of  $x$ . If  $x$  is a local minimizer, then it must satisfy the first- and second-order *necessary* optimality conditions. Conversely, a point  $x$  satisfying only the first-order condition is called a *first-order critical point*, while a point satisfying both of the conditions is called a *second-order critical point*. We also call it a *solution*. We call a solution *spurious* if it is not a global minimum. In this work, we study how existence of a spurious solution depends on the size/volume of the domain as well as the underlying structure of the problem.

The analysis of the landscape of the objective function around a global optimum may lead to an optimality guarantee for local search algorithms initialized sufficiently close to the solution (Keshavan, Montanari, and Oh 2010a, 2010b; Jain, Netrapalli, and Sanghavi 2013; Zheng and Lafferty 2015; Zhao, Wang, and Liu 2015; Sun and Luo 2016). Finding a good initialization scheme is highly problem-specific and difficult to generalize. Global analysis of the landscape is harder, but potentially more rewarding.

Both local and global convergence guarantees have been developed to justify the success of local search methods in various applications like dictionary learning (Agarwal et al. 2016), basic non-convex M-estimators (Mei, Bai, and Montanari 2016), shallow (Soltanolkotabi 2017) and deep (Yun, Sra, and Jadbabaie 2018) artificial neural networks with different activation (Li, Ding, and Sun 2018) and loss (Nouiehed and Razaviyayn 2018) functions, phase retrieval (Chen et al. 2018; Vaswani, Nayer, and Eldar 2017; Candes, Li, and Soltanolkotabi 2015) and more

general matrix sensing problems (Ge, Jin, and Zheng 2017; Josz et al. 2018). Particularly, significant progress has been made towards understanding different variants of *low-rank matrix recovery*, although explanations of the simplest version called *matrix sensing* are still under active development (Zhu et al. 2018; Chen et al. 2019; Li et al. 2019; Ge, Jin, and Zheng 2017; Chi, Lu, and Chen 2018; Zhang et al. 2018). Given a linear sensing operator  $\mathcal{A} : \mathbb{S}^n \rightarrow \mathbb{R}^m$  and a ground truth matrix  $z \in \mathbb{R}^{n \times r}$  ( $r < n$ ), an instance of the rank- $r$  matrix sensing problem consists in minimizing over  $\mathbb{R}^{n \times r}$  the nonconvex function

$$f_{z, \mathcal{A}}(x) = \|\mathcal{A}(xx^T - zz^T)\|_2^2 = \|\mathcal{A}(xx^T) - b\|_2^2, \quad (1)$$

where  $b = \mathcal{A}(zz^T)$ . We consider this function over a general set  $\mathcal{X} \subseteq \mathbb{R}^{n \times r}$ , although in this section we set  $\mathcal{X} = \mathbb{R}^{n \times r}$ . Recent work has generally found a certain assumption on the sensing operator to be sufficient for the matrix sensing problem to be “computationally easy to solve”. Precisely, this assumption works with the notion of RIP.

**Definition 1 (Restricted Isometry Property)** The linear map  $\mathcal{A} : \mathbb{S}^n \rightarrow \mathbb{R}^m$  is said to satisfy  $\delta_r$ -RIP for some constant  $\delta_r \in [0, 1)$  if there is  $\gamma > 0$  such that

$$(1 - \delta_r)\|X\|_F^2 \leq \gamma\|\mathcal{A}(X)\|_2^2 \leq (1 + \delta_r)\|X\|_F^2$$

holds for all  $X \in \mathbb{S}^n$  satisfying  $\text{rank}(X) \leq r$ .

The existing results proving absence of spurious local minima using this notion (such as Ge et al. 2015; Sun, Qu, and Wright 2015, 2018; Bhojanapalli, Neyshabur, and Srebro 2016; Ge, Lee, and Ma 2016; Ge, Jin, and Zheng 2017; Park et al. 2016; Zhu et al. 2018) are based on a norm-preserving argument: the problem turns out to be a low-dimensional embedding of a canonical problem known to contain no spurious local minima. While the approach is widely applicable in its scope, it requires fairly strong assumptions on the data. In contrast, Zhang et al. 2018; Zhang, Sojoudi, and Lavaei 2019 introduced a technique to find a certificate to guarantee that any given point cannot be a spurious local minimum of the problem of minimizing  $f_{z, \mathcal{A}}$  over the set  $\mathcal{X} \subseteq \mathbb{R}^{n \times r}$ , where  $z \in \mathcal{Z} \subseteq \mathbb{R}^{n \times r}$  and  $\mathcal{A}$  satisfies  $\delta_{2r}$ -RIP. Note that two different sets  $\mathcal{X}$  and  $\mathcal{Z}$  are involved here. Since  $f_{z, \mathcal{A}}$  depends on  $z$  and  $\mathcal{A}$ , this introduces a class of optimization problems defined as

$$\left\{ \min_{x \in \mathcal{X}} f_{z, \mathcal{A}}(x) \mid \mathcal{A} \text{ satisfies } \delta_{2r}\text{-RIP}, z \in \mathcal{Z} \right\}. \quad (\text{Problem}^{\text{RIP}})$$

(Problem<sup>RIP</sup>) consists of infinitely many instances of an optimization problem, each corresponding to some point  $z$  in  $\mathcal{Z}$  and some operator  $\mathcal{A}$  satisfying  $\delta_{2r}$ -RIP. The state-of-the-art results for (Problem<sup>RIP</sup>) are stated below.

**Theorem 1 (Bhojanapalli, Neyshabur, and Srebro 2016; Ge, Jin, and Zheng 2017; Zhang, Sojoudi, and Lavaei 2019)** *By taking  $\mathcal{X} = \mathcal{Z} = \mathbb{R}^{n \times r}$ , the following statements hold:*

- *If  $\delta_{2r} < 1/5$ , no instance of (Problem<sup>RIP</sup>) has a spurious second-order critical point.*

- If  $r = 1$  and  $\delta_2 < 1/2$ , then no instance of (Problem<sup>RIP</sup>) has a spurious second-order critical point.
- If  $r = 1$  and  $\delta_2 \geq 1/2$ , then there exists an instance of (Problem<sup>RIP</sup>) with a spurious second-order critical point.

Non-existence of a spurious second-order critical point effectively means that any algorithm that converges to a second-order critical point is guaranteed to recover  $zz^T$  exactly. Examples of such algorithms include variants of the stochastic gradient descent (SGD) that is known to avoid saddle or even spurious local minimum points under certain assumptions (Daneshmand et al. 2018), and widely used in machine learning (Krizhevsky, Sutskever, and Hinton 2012; Bottou and Bousquet 2008). Besides SGD, many local search methods have been shown to be convergent to second-order critical points with high probability under mild conditions, including the classical gradient descent (Lee et al. 2016), alternating minimizations (Li, Zhu, and Tang 2019) and Newton’s method (Paternain, Mokhtari, and Ribeiro 2019). In this paper, we present guarantees on the global optimality of the second-order critical points, which means that our results can be combined with any of the algorithms mentioned above to guarantee the global convergence.

Theorem 1 discloses the limits on the guarantees that the notion of RIP can provide. However, linear maps in applications related to physical systems, such as power system analysis, typically have no RIP constant smaller than 0.9, and yet the non-convex matrix sensing still manages to work on those instances. This gap between theory and practice motivates the following question.

**What is the alternative property practical problems satisfy that makes them easy to solve via simple local search?**

This question was studied earlier for special cases of matrix sensing, namely phase retrieval (Sun, Qu, and Wright 2018) and matrix completion (Ge, Lee, and Ma 2016). In case of phase retrieval, the alternative property consists in the particular distribution of the measurements operator. In matrix completion, the assumption includes conditions on the properties of the matrix being recovered along with conditions on the measurement operator itself. We address this question by developing a new notion that deals with the measurement operator and precisely captures when a structured matrix recovery problem has no spurious solution over an arbitrary ball. We focus the analysis over a given ball since local search methods tend to search over a neighborhood rather than the entire space, based on prior knowledge. In Section 2, we motivate the need for a new notion replacing or improving RIP with real-world examples. Section 3 introduces some formal definitions and develops a mathematical framework to analyze spurious solutions and relate them to the underlying sparsity and structure of the problem, using techniques in conic optimization. Sections 4 and 5 give the theory behind this notion and examples of its application. In Section 6, we present numerical results of the application of the developed theory to a real-world problem appearing in power systems analysis. Concluding remarks are given in Section 7. Some of the proofs, technical details and lemmas are collected in the Appendix.

## Notation

$\mathbb{C}^n$ ,  $\mathbb{R}^n$  and  $\mathbb{R}^{n \times r}$  denote the sets of complex and real  $n$ -dimensional vectors, and  $n \times r$  matrices, respectively.  $\mathbb{S}^n$  denotes the set of  $n \times n$  symmetric matrices.  $\text{Tr}(A)$ ,  $\|A\|_F$  and  $\langle A, B \rangle$  are the trace of a square matrix  $A$ , its Frobenius norm, and the Frobenius inner product of matrices  $A$  and  $B$  of compatible sizes. The normal distribution with mean  $\mu$  and covariance matrix  $\Sigma$  is denoted as  $\mathcal{N}(\mu, \Sigma)$ . In any linear space,  $\mathbf{1}$  is a vector whose entries are all equal to 1 and  $I$  is the identity matrix. For  $\omega \in \mathbb{R}^{n \times r}$  and  $R \in \mathbb{R} \cup \{+\infty\}$ , we define  $\mathbb{B}_R(\omega) = \{a \in \mathbb{R}^{n \times r} : \|a - \omega\|_F \leq R\}$ ,  $\bar{\mathbb{B}}_R(\omega) = \{a \in \mathbb{R}^{n \times r} : \|a - \omega\|_F < R\}$  and  $\partial \mathbb{B}_R(\omega) = \{a \in \mathbb{R}^{n \times r} : \|a - \omega\|_F = R\}$ . It follows from the definition that  $\partial \mathbb{B}_{+\infty}(\omega) = \emptyset$  and minimization over this set results in  $+\infty$  for any objective function. For a square matrix  $A$ , we define the symmetric part  $\text{Sym}(A) = (A + A^T)/2$ . For a symmetric matrix  $A$ , its null space is denoted with  $\text{Ker}(A)$ . For square matrices  $A_1, A_2, \dots, A_n$ , the matrix  $\text{diag}(A_1, \dots, A_n)$  is block-diagonal, with  $A_i$ 's on the block diagonal. The notation  $A \circ B$  refers to the Hadamard (entrywise) multiplication, and  $A \otimes B$  refers to the Kronecker product of matrices. The vectorization operator  $\text{vec} : \mathbb{R}^{n \times r} \rightarrow \mathbb{R}^{nr}$  stacks the columns of a matrix into a vector. The *matricization* operator  $\text{mat}(\cdot)$  is the inverse of  $\text{vec}(\cdot)$ . Let  $\succeq$  denote the positive semidefinite sign.

For a linear operator  $\mathcal{L} : \mathbb{R}^{n \times r} \rightarrow \mathbb{R}^m$ , the adjoint operator is denoted by  $\mathcal{L}^T : \mathbb{R}^m \rightarrow \mathbb{R}^{n \times r}$ . The matrix  $\mathbf{L} \in \mathbb{R}^{m \times nr}$  such that  $\mathcal{L}(x) = \mathbf{L} \text{vec}(x)$  is called the *matrix representation* of the linear operator  $\mathcal{L}$ . Bold letters are reserved for matrix representations of corresponding linear operators.

*Sparsity pattern*  $S$  of a set of matrices  $\mathbf{M} \subset \mathbb{R}^{m \times n}$  is a subset of  $\{1, \dots, \max\{n, m\}\}^2$  such that  $(i, j) \in S$  if and only if there is  $X \in \mathbf{M}$  with the property that  $X_{ij} \neq 0$ . Given a sparsity pattern  $S$ , define its matrix representation  $\mathbf{S} \in \mathbb{S}^{m \times n}$  as

$$S_{ij} = \begin{cases} 0 & \text{if } (i, j) \in S, \\ 1 & \text{if } (i, j) \notin S, \end{cases}$$

The *orthogonal basis* of a given  $m \times n$  matrix  $A$  (with  $m \geq n$ ) is a matrix  $P = \text{orth}(A) \in \mathbb{R}^{m \times \text{rank}(A)}$  consisting of  $\text{rank}(A)$  orthonormal columns that span  $\text{range}(A)$ :

$$P = \text{orth}(A) \iff PP^T A = A, P^T P = I_{\text{rank}(A)}.$$

Positive part means  $(\cdot)_+ = \max\{0, \cdot\}$ , and eigenvalues in an arbitrary order are denoted by  $\lambda_i(\cdot)$ .

## 2 Motivating example

In this section, we motivate this work by offering a case study on data analytics for energy systems. The state of a power system can be modeled by a vector of complex voltages on the nodes (buses) of the network. Monitoring the state of a power system is obviously a necessary requirement for its efficient and safe operation. This crucial information should be inferred from some measurable parameters, such as the power that is generated and consumed at each bus or transmitted through a line. The power network can be modeled by a number of parameters grouped into the admittance

matrix  $Y \in \mathbb{C}^{n \times n}$ . The state estimation problem consists in recovering the unknown voltage vector  $v \in \mathbb{C}^n$  from the available measurements. In the noiseless scenario, these measurements are  $m$  real numbers of the form

$$v^* M_i v, \quad \forall i \in \{1, \dots, m\}, \quad (2)$$

where  $M_i = M_i(Y) \in \mathbb{C}^{n \times n}$  are sparse Hermitian matrices representing power-flow and power-injection as well as voltage magnitudes measurements. The sparsity pattern of the measurement matrices is determined by the topology of the network, while its nonzero entries are certain known functions of the entries of  $Y$ . Since the total number of nonzero elements in matrices  $M_i$  exceeds the total number of parameters contained in  $Y$ , we can think of  $Y \rightarrow \{M_i\}_{i=1}^m$  as an embedding from a low-dimensional space. For a detailed discussion on the problem formulation and approaches to its solution see e.g. Zhang, Madani, and Lavaei 2018.

To formulate the problem as a low-rank matrix recovery, we introduce a sparse matrix  $\mathbf{A} = \mathbf{A}(Y) \in \mathbb{C}^{m \times n^2}$  with  $i$ -th row equal to  $\text{vec}(M_i)^T$ . The measurement vector can be written as  $\mathbf{A} \text{vec}(vv^T)$ . To find  $v$  from the measurements, one can solve the non-convex optimization problem:

$$\underset{x \in \mathbb{C}^n, \|x - \omega\|_F \leq R}{\text{minimize}} \quad \|\mathbf{A} \text{vec}(xx^T - vv^T)\|_F^2 \quad (3)$$

where  $\omega \in \mathbb{C}^n$  and  $R \in \mathbb{R} \cup \{+\infty\}$  are some parameters determined by the prior knowledge about the solution  $v$ . In practice, this non-convex optimization problem is usually solved via local search methods, which converge to a second-order critical point at best. Since  $f(x) = \|\mathbf{A} \text{vec}(xx^T - vv^T)\|_F^2 = \langle xx^T - vv^T, \mathbf{A}^T \mathbf{A} \text{vec}(xx^T - vv^T) \rangle$ , the set of critical points for the problem is defined by the linear map represented with the matrix  $\mathbf{H} = \mathbf{A}^T \mathbf{A}$ , which thus is the key subject of the study. Problems arising in power systems analysis are based on operators that possess a specific structure. An example of a structure for the matrix  $\mathbf{A}$  is given in Fig. 1a, and the structure of the corresponding  $\mathbf{H}$  is described in Fig. 1b. The respective power network will be considered in more details in Section 6. As discussed previously, given  $\mathbf{H}$ , it is practically important to know if there exist  $v, x \in \mathbb{C}^n$  such that  $x$  is a critical point of (3) while  $xx^T \neq vv^T$ . Absence of these points proves that a local search method recovers  $v$  exactly, certifying safety of its use. For example, in case of unconstrained optimization ( $R = +\infty$ ), the answer can be provided by the following problem having its optimal objective value equal to zero:

$$\begin{aligned} & \underset{v, x \in \mathbb{C}^n}{\text{minimize}} \quad \|\mathcal{A}(xx^T - vv^T)\|_F^2 \\ & \text{subject to} \quad \nabla_x f(x) = 0 \\ & \quad \quad \quad \nabla_x^2 f(x) \succeq 0 \end{aligned}$$

However, this is an  $\mathcal{NP}$ -hard problem in general and cannot be solved efficiently. Even if we solved it, the sensing operator  $\mathcal{A}$  could change over time without changing its structure, and therefore any conclusion made for a specific problem cannot be generalized to other ones that should be solved for real-world problems where data analysis is to be performed periodically. One way to circumvent this issue is to develop a sufficient condition for all mapping  $\mathbf{H}$  with the same structure.

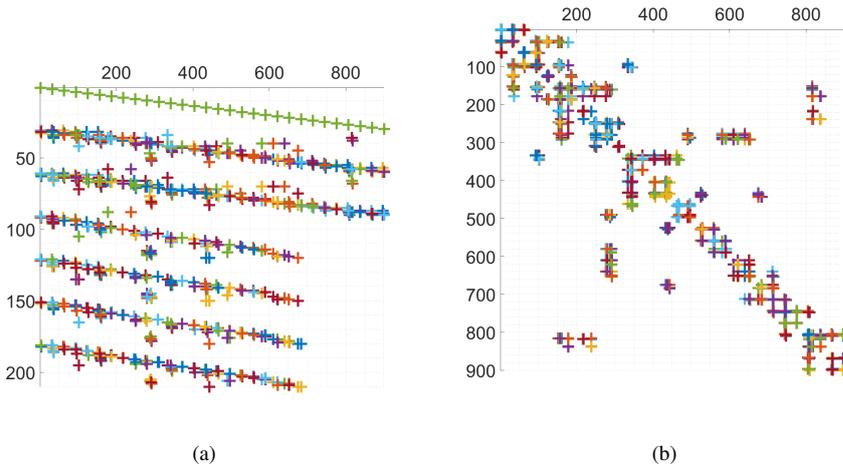


Fig. 1: Examples of the structure patterns of operators  $\mathcal{A}$  (left plot) and  $\mathcal{H}$  (right plot) in power system applications. The positions of the identical nonzero entries of a matrix are marked with the same markers.

### 3 Introducing Kernel Structure

Consider a linear operator  $\mathcal{A} : \mathbb{S}^n \rightarrow \mathbb{R}^m$  with the matrix representation  $\mathbf{A} \in \mathbb{R}^{m \times n^2}$  and a sparsity pattern  $S_{\mathcal{A}}$ . Assume that there is a set of hidden parameters  $\xi \in \mathbb{R}^d, d \ll m$ , such that  $\mathbf{A}$  is the image of  $\omega$  in the space of a much higher dimension. In this way,  $\mathcal{A}$  has a low-dimensional structure beyond sparsity, which is captured by  $\mathbf{A} = \mathbf{A}(\xi)$  and  $\mathbf{A}(0) = \mathbf{0}$ . The motivating example in Section 2 is a special case of this construction since it could be stated entirely with the real vectors and matrices of a bigger size. We define the nonconvex objective

$$f : \mathbb{R}^{n \times r} \rightarrow \mathbb{R} \quad \text{such that} \quad f(x) = \|\mathcal{A}(xx^T - zz^T)\|^2$$

parametrized by  $\mathcal{A}$  and  $z \in \mathbb{R}^{n \times r}$ . Its value is always nonnegative by construction, and the global minimum 0 is attainable. To emphasize the dependence on certain parameters, we will write them in the subscript. To align the minimization problem with the problem of reconstructing  $zz^T$ , we need to introduce a regularity assumption:

**Assumption 1.** *The  $2r$ -RIP constant  $\delta_{2r}$  of  $\mathcal{A}$  exists (and by definition is strictly smaller than 1).*

Note that we do not assume any particular value for the RIP constant here. We will rely on Assumption 1 throughout the paper. This assumption implies that for all  $x, z \in \mathbb{R}^{n \times r}$ :

$$\|\mathcal{A}(xx^T - zz^T)\| = 0 \text{ if and only if } xx^T = zz^T$$

Another way to express the objective is

$$f(x) = \langle xx^T - zz^T, \mathcal{H}(xx^T - zz^T) \rangle.$$

Here,  $\mathcal{H} = \mathcal{A}^T \mathcal{A}$  is the linear *kernel operator* that has the matrix representation  $\mathbf{H} = \mathbf{A}^T \mathbf{A}$  and sparsity pattern  $S_{\mathcal{H}}$ . Namely,  $(i, j) \in S_{\mathcal{H}}$  if and only if there exists  $k$  such that  $(k, i) \in S_{\mathcal{A}}$  and  $(k, j) \in S_{\mathcal{A}}$ . Sparsity of  $\mathcal{H}$  is controlled by the out-degree of the graph represented by  $S_{\mathcal{A}}$ , and tends to be low in applications like power systems.  $S_{\mathcal{H}}$  is represented by a matrix  $\mathbf{S}$ , so that  $S_{\mathcal{H}}$ -sparse operators are exclusively those satisfying the linear equation  $\mathcal{S}(\mathbf{H}) = \mathbf{S} \circ \mathbf{H} = \mathbf{0}$ . Besides sparsity,  $\mathcal{H}$  inherits the low-dimensional structure from  $\mathcal{A}$ , which can be captured by  $\mathbf{H} = \mathbf{A}(\xi)^T \mathbf{A}(\xi) = \mathbf{H}(\xi)$  where  $\xi \in \mathbb{R}^d$ . This dependence can be locally approximated in the hidden parameter space with a linear one. More precisely, suppose that there is a linear operator  $\mathcal{W}$  defined over  $\mathbb{S}^{n^2}$  such that  $\mathcal{W}(\mathbf{H}(\xi)) \approx 0$  for the values of  $\xi$  under consideration. Thus, from now on we focus exclusively on low-dimensional structures of the form  $\mathcal{W}(\mathbf{H}) = 0$ . Together, the sparsity operator  $\mathcal{S}$  and the low-dimensional structure operator  $\mathcal{W}$  form the combined structure operator  $\mathcal{T} = (\mathcal{S}, \mathcal{W})$  that accumulates the structure of the kernel operator.

**Definition 2 (Kernel Structure Property or KSP)** The linear map  $\mathcal{A} : \mathbb{S}^n \rightarrow \mathbb{R}^m$  is said to satisfy  $\mathcal{T}$ -KSP if it satisfies Assumption 1 and there is a linear structure operator  $\mathcal{T} : \mathbb{S}^{n^2} \rightarrow \mathbb{R}^l$  such that

$$\mathcal{T}(\mathbf{A}^T \mathbf{A}) = 0$$

where  $\mathbf{A}$  is the matrix representation of  $\mathcal{A}$ .

Notice that a particular sensing operator  $\mathcal{A}$  can be kernel structured with respect to an entire family of structure operators, and we can possibly select any of them for our benefit in the following section.

#### 4 Using KSP

After fixing the kernel structure of the sensing operators,  $\omega \in \mathbb{R}^{n \times r}$  and  $R \in \mathbb{R} \cup \{+\infty\}$ , we can state the problem under study as follows:

$$\left\{ \min_{x \in \mathbb{B}_R(\omega)} f_{z, \mathcal{A}}(x) \mid \mathcal{A} \text{ satisfies Assumption 1 and } \mathcal{T}\text{-KSP}, z \in \mathbb{B}_R(\omega) \right\},$$

(Problem<sup>KSP</sup>)

Note that (Problem<sup>KSP</sup>) consists of infinitely many instances of an optimization problem, each corresponding to some point  $z \in \mathbb{B}_R(\omega)$  and some operator  $\mathcal{A}$  satisfying  $\mathcal{T}$ -KSP.

If  $X$  is regarded as an input and the operator  $\mathcal{A}$  is regarded as a system with its output being  $\mathcal{A}(X)$ , the RIP constant aims at characterizing the input-output behavior of the system. This input-output relationship can also be controlled by imposing the following constraint on the matrix  $\mathcal{H}$ :

$$(1 - \delta)\mathcal{I} \preceq \mathcal{H} \preceq (1 + \delta)\mathcal{I},$$

where  $\mathcal{I}$  is the identity operator. More precisely, the above inequality guarantees that the operator  $\mathcal{A}$  has an RIP constant less than or equal to  $\delta$ . Inspired by this

observation, we introduce the function  $\mathbb{O}(x, z; \mathcal{T})$  to be the optimal objective value of the convex optimization problem:

$$\text{minimum}_{\delta \in \mathbb{R}, \mathcal{H}} \delta$$

$$\text{subject to } \mathcal{L}_{x,z}(\mathcal{H}) = 0 \quad (4a)$$

$$\mathcal{M}_{x,z}(\mathcal{H}) \succeq 0 \quad (4b)$$

$$\mathcal{T}(\mathcal{H}) = 0 \quad (4c)$$

$$(1 - \delta)\mathcal{I} \preceq \mathcal{H} \preceq (1 + \delta)\mathcal{I} \quad (4d)$$

where  $\mathcal{L}_{x,z}(\mathcal{H}) = \nabla f_{z,\mathcal{H}}(x)$  and  $\mathcal{M}_{x,z}(\mathcal{H}) = \nabla^2 f_{z,\mathcal{H}}(x)$ . This optimization is performed over all operators  $\mathcal{H}$  satisfying the KSP. We will later show that the function  $\mathbb{O}$  sets an upper bound on the  $\delta_{2r}$  such that none of the functions  $f_{z,\mathcal{A}}$  with  $\mathcal{A}$  satisfying  $\mathcal{T}$ -KSP and  $\delta_{2r}$ -RIP has a spurious second-order critical point at  $x$ , provided that  $x$  is not on the boundary of the optimization domain  $\mathbb{B}_R(\omega)$ .

For completeness, and for further reference within this paper, we calculate the analytic forms of the first- and second-order derivatives of  $f_{z,\mathcal{H}}$  below. Introduce a vector  $\mathbf{e}$  and a matrix  $\mathbf{X}$  such that for all  $u \in \mathbb{R}^{n \times r}$  it holds that

$$\mathbf{e} = \text{vec}(xx^T - zz^T), \quad \mathbf{X}\text{vec}(u) = \text{vec}(xu^T + ux^T).$$

We write the operators  $\mathcal{L}$ ,  $\mathcal{M}$  and their transpose operators in closed form:

$$\begin{aligned} \mathcal{L}_{x,z} : \mathbb{S}^{n^2} &\rightarrow \mathbb{R}^{n \times r} & \mathcal{L}_{x,z}(\mathbf{H}) &= 2 \cdot \mathbf{X}^T \mathbf{H} \mathbf{e}, \\ \mathcal{L}_{x,z}^T : \mathbb{R}^{n \times r} &\rightarrow \mathbb{S}^{n^2} & \mathcal{L}_{x,z}^T(y) &= \mathbf{e}y^T \mathbf{X}^T + \mathbf{X}y\mathbf{e}^T, \\ \mathcal{M}_{x,z} : \mathbb{S}^{n^2} &\rightarrow \mathbb{S}^{nr} & \mathcal{M}_{x,z}(\mathbf{H}) &= [I_r \otimes (\text{mat}(\mathbf{H}\mathbf{e}) + \text{mat}(\mathbf{H}\mathbf{e})^T)] + \mathbf{X}^T \mathbf{H} \mathbf{X}, \\ \mathcal{M}_{x,z}^T : \mathbb{S}^{nr} &\rightarrow \mathbb{S}^{n^2} & \mathcal{M}_{x,z}^T(V) &= \text{vec}(V)\mathbf{e}^T + \text{evec}(V)^T + \mathbf{X}V\mathbf{X}^T. \end{aligned}$$

Since  $f_{z,\mathcal{H}}(x)$  is linear in  $\mathcal{H}$ , the operators  $\mathcal{L}_{x,z}$  and  $\mathcal{M}_{x,z}$  are both linear. Thus, the problem defining the function  $\mathbb{O}$  is convex.

The function  $\mathbb{O}$  is useful only for those points  $x$  that are located strictly inside the domain  $\mathbb{B}_R(\omega)$ . This is due to the fact the constraints (4a) and (4b) are meant to be optimality conditions for a point inside the domain  $\mathbb{B}_R(\omega)$ . For a point  $x$  that lies on the boundary, we define the corresponding function  $\mathbb{O}^{\partial\mathbb{B}}(x, z; \mathcal{T}, \omega)$  as the optimal objective value of the convex optimization problem:

$$\text{minimum}_{\delta, \mu \geq 0, \mathcal{H}} \delta$$

$$\text{subject to } \mathcal{L}_{x,z}(\mathcal{H}) = -\mu(x - \omega) \quad (5a)$$

$$P_{x-\omega} \mathcal{M}_{x,z}(\mathcal{H}) P_{x-\omega}^T \succeq 0 \quad (5b)$$

$$\mathcal{T}(\mathcal{H}) = 0 \quad (5c)$$

$$(1 - \delta)\mathcal{I} \preceq \mathcal{H} \preceq (1 + \delta)\mathcal{I} \quad (5d)$$

where  $P_{x-\omega} \in \mathbb{R}^{(nr-1) \times nr}$  is the matrix of orthogonal projection onto the subspace orthogonal to  $x - \omega$ . The role of  $\mathbb{O}^{\partial\mathbb{B}}$  is the same as  $\mathbb{O}$  but will be used only for those

values of  $x$  such that  $\|x - \omega\| = R$ . Note that (5a) and (5b) are the necessary optimal conditions for a point on the boundary of  $\mathbb{B}_R(\omega)$ .

To relax the  $\delta_{2r}$ -RIP condition, we consider those operators that have a bounded effect on a linear subspace of limited-rank inputs. Indeed, for any  $2r$  linearly independent vectors, the linear span of them is a linear subspace of the manifold of the  $2r$ -rank matrices. Thus, for any linear operator  $\mathcal{P}$  from a  $2r$ -dimensional (or lower) vector space to  $\mathbb{R}^{n^2}$ , the following condition on  $\mathcal{H}$  holds if  $\mathcal{A}$  satisfies  $\delta$ -RIP:

$$(1 - \delta) \mathcal{P}^T \mathcal{P} \preceq \mathcal{P}^T \mathcal{H} \mathcal{P} \preceq (1 + \delta) \mathcal{P}^T \mathcal{P}. \quad (6)$$

Based on this observation, we define the function  $\mathbb{O}_P(x, z; \mathcal{T})$  as the optimal objective value of the following convex optimization problem:

$$\begin{aligned} & \text{minimum}_{\delta \in \mathbb{R}, \mathcal{H}} \quad \delta \end{aligned} \quad (7a)$$

$$\text{subject to } \mathcal{L}_{x,z}(\mathcal{H}) = 0 \quad (7a)$$

$$\mathcal{M}_{x,z}(\mathcal{H}) \succeq 0 \quad (7b)$$

$$\mathcal{T}(\mathcal{H}) = 0 \quad (7c)$$

$$(1 - \delta) \mathcal{P}^T \mathcal{P} \preceq \mathcal{P}^T \mathcal{H} \mathcal{P} \preceq (1 + \delta) \mathcal{P}^T \mathcal{P} \quad (7d)$$

where  $\mathcal{P}$  is the linear operator from  $\mathbb{R}^{\text{rank}([x \ z])^2}$  to  $\mathbb{R}^{n^2}$  that is represented by the matrix  $\mathbf{P} = \text{orth}([x \ z]) \otimes \text{orth}([x \ z])$ . Note that (7) is obtained from (4) by replacing its constraint (4d) with the milder condition (6). We will show that the function  $\mathbb{O}_P$  sets a lower bound on the  $\delta_{2r}$  such that none of the functions  $f_{z;\mathcal{A}}$  with  $\mathcal{A}$  satisfying  $\mathcal{T}$ -KSP and  $\delta_{2r}$ -RIP has a spurious second-order critical point at  $x$ , provided that  $x$  is not on the boundary of the optimization domain  $\mathbb{B}_R(\omega)$ .

Similarly to  $\mathbb{O}^{\partial \mathbb{B}}$ , the function  $\mathbb{O}_P^{\partial \mathbb{B}}(x, z; \mathcal{T}, \omega)$  is defined as the optimal objective value of the convex optimization problem:

$$\begin{aligned} & \text{minimum}_{\delta, \mu > 0, \mathcal{H}} \quad \delta \\ & \text{subject to } \mathcal{L}_{x,z}(\mathcal{H}) = -\mu(x - \omega) \\ & \quad P_{x-\omega} \mathcal{M}_{x,z}(\mathcal{H}) P_{x-\omega}^T \succeq 0 \\ & \quad \mathcal{T}(\mathcal{H}) = 0 \\ & \quad (1 - \delta) \mathcal{P}^T \mathcal{P} \preceq \mathcal{P}^T \mathcal{H} \mathcal{P} \preceq (1 + \delta) \mathcal{P}^T \mathcal{P} \end{aligned}$$

which is designed to lower bound the  $\delta_{2r}$  such that none of the functions  $f_{z;\mathcal{A}}$  with  $\mathcal{A}$  satisfying  $\mathcal{T}$ -KSP and  $\delta_{2r}$ -RIP has a spurious second-order critical point at  $x$ , provided that  $x$  is on the boundary of  $\mathbb{B}_R(\omega)$ .

Now, we are ready to state one of the main results of this paper.

**Theorem 2 (KSP necessary and sufficient conditions)** *For all instances of (Problem<sup>KSP</sup>), there are no spurious second-order critical points if*

$$\begin{cases} \mathbb{O}_P(x, z; \mathcal{T}) \equiv 1 \text{ over } \mathbb{B}_R(\omega) \times \mathbb{B}_R(\omega) \setminus \{xx^T = zz^T\} \\ \mathbb{O}_P^{\partial \mathbb{B}}(x, z; \mathcal{T}, \omega) \equiv 1 \text{ over } \partial \mathbb{B}_R(\omega) \times \mathbb{B}_R(\omega) \setminus \{xx^T = zz^T\} \end{cases} \quad (8)$$

and only if

$$\begin{cases} \mathbb{O}(x, z; \mathcal{T}) \equiv 1 \text{ over } \mathbb{B}_R(\omega) \times \mathbb{B}_R(\omega) \setminus \{xx^T = zz^T\} \\ \mathbb{O}^{\partial \mathbb{B}}(x, z; \mathcal{T}, \omega) \equiv 1 \text{ over } \partial \mathbb{B}_R(\omega) \times \mathbb{B}_R(\omega) \setminus \{xx^T = zz^T\} \end{cases} \quad (9)$$

This theorem is formally proven in the Appendix. To elaborate on implications and practicality of the result, we present its application for a specific structure of the sensing operator below.

#### 4.1 Ellipsoid norm: Rank 1

In this subsection, we prove a special case of Theorem 2 for the ellipsoid norm objective function and  $R = +\infty$ . This proof first provides useful intuition behind the proof of the general case and then simplifies the conditions of Theorem 2 to show that they always hold for a specific class of operators.

Consider the ellipsoid norm of  $xx^T - zz^T$  given by a full-rank matrix  $Q \in \mathbb{R}^{n \times n}$ , denoted with  $h$ :

$$h(x) = \|Q(xx^T - zz^T)\|_F^2 = f_{z, \mathbf{A}}(x)$$

With no loss of generality, assume that  $Q \in \mathbb{S}^n$  since  $h(\cdot)$  really depends only on  $Q^T Q$ . The function can be implemented with a block-diagonal sensing operator matrix  $\mathbf{A} = \text{diag}(Q, \dots, Q) \in \mathbb{S}^{n^2}$ , which generates a block-diagonal kernel matrix  $\mathbf{H} = \text{diag}(QQ, \dots, QQ)$ . Thus, the kernel matrix is a block-diagonal matrix  $\mathbf{H} = \text{diag}(H_{11}, \dots, H_{mm}) \in \mathbb{S}^{n^2}$  with blocks of size  $n \times n$  equal to each other; in other words,  $H_{ii} = H_{jj}$  for all  $i, j \in \{1, \dots, n\}$ . This generates a kernel structure. By applying the theory introduced above, we obtain the following result for the rank-one case.

**Proposition 1** Consider a kernel structure operator  $\mathcal{T} = (\mathcal{S}, \mathcal{W})$  such that

- $\mathcal{S}(\mathbf{H}) = \mathbf{0}$  iff  $\mathbf{H} = \text{diag}(H_{11}, \dots, H_{mm})$
- $\mathcal{W}(\mathbf{H}) = \mathbf{0}$  iff  $H_{ii} = H_{jj}, i, j \in \{1, \dots, n\}$ ,

Then, no instance of the (Problem<sup>KSP</sup>) with  $R = \infty$  has a spurious second-order critical point over  $\mathbb{R}^n$ .

The proposition implies that the function  $h(x)$  can never have a spurious solution for rank-1 arguments. Indeed, Assumption 1 and the following lemma combined imply that  $H_{ii}$ , and, consequently, its decomposition  $H_{ii} = QQ$ , are full-rank matrices.

**Lemma 1** For any  $\delta_r \in [0, 1)$ , the matrix  $Q \in \mathbb{S}^n$  satisfies

$$(1 - \delta_r) \|X\|_F^2 \leq \|QX\|_F^2 \leq (1 + \delta_r) \|X\|_F^2$$

for every  $X$  such that  $\text{rank}(X) \leq r$  only if  $\text{rank}(Q) = n$

*Proof.* For contradiction, suppose that  $u \in \text{Ker}(Q)$  and  $u \neq 0$ . Take  $X = uu^T$  and observe that  $QX = 0$ , which contradicts that  $(1 - \delta_r) \|X\|_F^2 \leq \|QX\|_F^2$ .  $\square$

The following lemma provides a version the conditions (8) and (9) combined for this particular structure operator.

**Lemma 2** *Given  $z \in \mathbb{R}^{n \times r}$ , a point  $x \in \mathbb{R}^{n \times r}$  is not a first-order critical point of the function  $h(\cdot)$  for an arbitrary full-rank matrix  $Q$  if and only if there is  $\lambda \in \mathbb{R}^{n \times r}$  such that*

$$0 \neq \text{Sym}[(x\lambda^T + \lambda x^T)(xx^T - zz^T)] \succeq 0$$

*Proof.* By expanding  $h(x+u)$  as

$$h(x+u) = h(x) + \text{Tr}(2x^T((xx^T - zz^T)M + M(xx^T - zz^T))u) + \text{Tr}(u^T((xx^T - zz^T)M + M(xx^T - zz^T))u + (xu^T + ux^T)M(xu^T + ux^T)) + o(\|u\|^2)$$

one can arrive at a more specified expression for the second-order necessary conditions for local optimality:

$$\langle \nabla h(x), u \rangle = 2\langle Q(xx^T - zz^T), Q(xu^T + ux^T) \rangle = 0 \quad (10a)$$

$$\langle \nabla^2 h(x)u, u \rangle = 2\langle QQ(xx^T - zz^T), uu^T \rangle + \|Q(xu^T - ux^T)\|_F^2 \geq 0 \quad (10b)$$

for all  $u \in \mathbb{R}^{n \times r}$ . We re-arrange the first-order condition (10a):

$$\left( (xx^T - zz^T)M + M(xx^T - zz^T) \right) z = 0 \quad (11)$$

If the equation (11) does not hold for some  $M \succ 0$ , then  $x$  cannot be a critical point for that  $z$  and  $M$ . Consequently, the problem

$$\begin{aligned} & \text{minimize} && -\alpha \\ & M \in \mathbb{S}^n, \alpha \in \mathbb{R} \end{aligned}$$

$$\text{subject to} \quad \left( (yy^T - xx^T)M + M(yy^T - xx^T) \right) y = 0 \quad (12a)$$

$$M - \alpha I \succeq 0, \quad (12b)$$

is bounded by 0 if and only if the equation (11) does not hold for arbitrary  $M \succ 0$ . If  $x$  is critical for some  $M \succ 0$ , then it is unbounded.

The problem (12) is a semidefinite program with a zero duality gap, since  $M = 0$  and  $\alpha = -1$  constitute a strictly feasible primal point. We introduce the dual variable  $\lambda \in \mathbb{R}^{n \times r}$  for the equality constraint (12a) and the dual variable  $G \in \mathbb{S}^n$  for the Positive Semidefiniteness (PSD) constraint (12b). The dual problem can be written as

$$\max_{\lambda \in \mathbb{R}^{n \times r}, G \succeq 0} \min_{M \in \mathbb{S}^n, \alpha \in \mathbb{R}} \text{Tr}[(2\text{Sym}[(y\lambda^T + \lambda y^T)(yy^T - xx^T)] - G)M] + \alpha(\text{Tr}(G) - 1)$$

The inner optimization problem has a finite solution if and only if

$$\begin{cases} G &= (y\lambda^T + \lambda y^T)(yy^T - xx^T) + (yy^T - xx^T)(y\lambda^T + \lambda y^T) \\ \text{Tr}(G) &= 1 \end{cases}$$

The dual problem can be expressed as

$$\begin{aligned} & \text{maximize} && 0 \\ & \lambda \in \mathbb{R}^{n \times r}, G \in \mathbb{S}^n \end{aligned} \quad \begin{aligned} & \text{subject to} && G = \text{Sym}[(y\lambda^T + \lambda y^T)(yy^T - xx^T)], \\ & && \text{Tr}(G) = 1, \\ & && G \succeq 0 \end{aligned}$$

This is feasible if and only if the primal problem (12) is bounded. Consequently, it is feasible if and only if the point  $x \in \mathbb{R}^{n \times r}$  is not a critical point of the function  $h$  for all  $M \succ 0$ .

To eliminate the condition on trace, notice that a PSD matrix had a nonnegative trace, which is equal to zero if and only if the matrix is the zero matrix. All the constraints are homogeneous in  $G$ , so the trace can always be normalized to 1. Thus, the dual feasibility is equivalent to the condition  $0 \neq G \succeq 0$ . This concludes the proof.  $\square$

This condition will be relaxed further for simplicity below.

**Lemma 3** *Given  $z \in \mathbb{R}^{n \times r}$ , a point  $x \in \mathbb{R}^{n \times r}$  is not a first-order critical point of the function  $h(\cdot)$  for an arbitrary full-rank matrix  $Q$  if there are  $T_1 \in \mathbb{R}^{r \times r}$  and  $T_2 \in \mathbb{S}^r$  such that the matrix  $T = \begin{bmatrix} 0 & T_1 \\ -T_1^T & T_2 \end{bmatrix}$  satisfies the relations*

$$0 \neq [-z \ x] (T^T P + P T) \begin{bmatrix} -z^T \\ x^T \end{bmatrix} \succeq 0 \quad (13)$$

where  $P = \begin{bmatrix} z^T \\ x^T \end{bmatrix} [z \ x]$ .

*Proof.* Suppose that there exists  $T$  from the statement of the lemma. Notice that

$$\begin{aligned} xT_1^T z^T + \frac{1}{2} xT_2 x^T + zT_1 x^T + \frac{1}{2} xT_2 x^T &= [-z \ x] \begin{bmatrix} 0 & -T_1 \\ T_1^T & T_2 \end{bmatrix} \begin{bmatrix} z^T \\ x^T \end{bmatrix} = \\ &= [z \ x] \begin{bmatrix} 0 & T_1 \\ -T_1^T & T_2 \end{bmatrix} \begin{bmatrix} -z^T \\ x^T \end{bmatrix} \end{aligned}$$

and

$$xx^T - zz^T = [z \ x] \begin{bmatrix} -z^T \\ x^T \end{bmatrix} = [-z \ x] \begin{bmatrix} z^T \\ x^T \end{bmatrix}$$

We use the formulas above to expand the condition (13) and obtain

$$0 \neq \text{Sym}[(x(zT_1 + \frac{xT_2}{2})^T + (zT_1 + \frac{xT_2}{2})x^T)(xx^T - zz^T)] \succeq 0,$$

Conclusion immediately follows by applying Lemma 2 with  $\lambda = zT_1 + \frac{1}{2}xT_2$ .  $\square$

To prove Proposition 1, we check the condition above for all pairs of  $z$  and  $x$ .

Proof of Proposition 1.

We start by proving that any point except for 0 and  $\pm z$  cannot be a first-order critical point of the function  $h$ . Assume that  $x \notin \{0, \pm z\}$ . By Lemma 3, it is sufficient to prove that there are  $\alpha$  and  $\beta$  in  $\mathbb{R}$  such that the matrix  $T = \begin{bmatrix} 0 & \alpha \\ -\alpha & \beta \end{bmatrix}$  satisfies

$$0 \neq G = [-z \ x] (T^T P + PT) \begin{bmatrix} -z^T \\ x^T \end{bmatrix} \succeq 0$$

where  $P = \begin{bmatrix} z^T \\ x^T \end{bmatrix} [z \ x]$ . Consider three scenarios for  $x$  and  $z$ :

Case 1  $x = \gamma z$ :

$$G = [-z \ x] (T^T P + PT) \begin{bmatrix} -z^T \\ x^T \end{bmatrix} = 2\gamma(\gamma^2 - 1)(2\alpha + \beta\gamma)zz^T zz^T$$

For  $\alpha = \gamma(\gamma^2 - 1)$  and  $\beta = 0$ , it holds that  $G = (2\gamma(\gamma^2 - 1)zz^T)^2 \succeq 0$ . The matrix is nonzero for  $x \notin \{0, \pm z\}$ .

Case 2  $z^T x = 0$ :

The matrix  $P$  takes the form  $P = \begin{bmatrix} \|z\|^2 & 0 \\ 0 & \|x\|^2 \end{bmatrix}$ , so for  $\alpha = 0$  and  $\beta = 1$ , it holds that

$$G = 2\|x\|^2 xx^T \succeq 0$$

The matrix is nonzero for  $x \neq 0$ .

Case 3  $0 < (z^T x)^2 < \|z\|_2^2 \|x\|_2^2$ :

By scaling, we can assume without loss of generality that  $z^T x = 1$ ; thus  $P = \begin{bmatrix} \|z\|_2^2 & 1 \\ 1 & \|x\|_2^2 \end{bmatrix}$ .

It is sufficient to show that  $T^T P + PT \succ 0$  to guarantee  $G$  to be nonzero and PSD. To show this, we use Sylvester's criterion. The upper left corner of this matrix is equal to  $-2\alpha$ . Moreover,

$$\det(T^T P + PT) = -(\|x\|_2^2 - \|y\|_2^2)^2 - 4\alpha^2 - 2(\|x\|_2^2 + \|y\|_2^2)\alpha\beta - \beta^2$$

For  $\alpha = -1$ , the discriminant of this quadratic polynomial with respect to  $\beta$  is equal to  $D = 16(\|z\|_2^2 \|x\|_2^2 - 1)$ . By the strict Cauchy-Schwarz inequality in the assumption of the case,  $D$  is strictly greater than 0. Thus, there exists  $\beta$  such that the matrix is positive definite. This implies that none of  $x \notin \{0, \pm z\}$  satisfies the first-order necessary condition of local optimality for an unconstrained problem. Assume that  $x = 0$ . The quadratic form on the Hessian at this point

$$\langle \nabla^2 h(0)u, u \rangle = -2\langle Qz z^T, uu^T \rangle$$

takes a negative value at  $u = z$ . Thus, it does not satisfy the second-order necessary condition of local optimality for an unconstrained problem. The points  $x = \pm z$  are not spurious points, which concludes the proof.  $\square$

## 4.2 Ellipsoid norm: Higher ranks

The function  $h(\cdot)$  defined over  $\mathbb{R}^{n \times r}$  is significantly harder to study analytically if  $r > 1$ . Numerical application of Theorem 2 allows us to make a conjecture.

*Conjecture 1* For the kernel structure operator introduced in Proposition 1, no instance of the (Problem<sup>KSP</sup>) with  $\mathcal{L} = \mathbb{R}^{n \times r}$  has a spurious second-order critical points over  $\mathbb{R}^{n \times r}$  for an arbitrary  $r$ .

This conjecture is based on evaluation of  $\mathbb{O}(x, z; \mathcal{T})$  at 72000 pairs of points  $x, z \in \mathbb{R}^{8 \times 3}$  randomly sampled from the standard Gaussian distribution. All of them have the optimal value 1. However, if we consider first-order critical points as well, then it is straightforward to find a counterexample. After dropping the constraint on  $\mathcal{M}_{x,z}$  in the formulation of  $\mathbb{O}(x, z; \mathcal{T})$  and  $\mathbb{O}_P(x, z; \mathcal{T})$ , one can formulate a statement similar to Theorem 2 tailored to first-order solutions. The following proposition presents a corollary of the result.

**Proposition 2** For the kernel structure operator introduced in Proposition 1, for every  $n \geq 8$  and  $r > 1$ , there is  $z \in \mathbb{R}^{n \times r}$  such that (Problem<sup>KSP</sup>) has a spurious saddle point.

*Proof.* First, we prove it for  $n = 8$  and  $r = 2$  by a counterexample. Consider the two points

$$x = \begin{bmatrix} 0 & -1 \\ 1 & -1 \\ 1 & 1 \\ 1 & 0 \\ 0 & 0 \\ -1 & 1 \\ -1 & 1 \\ 0 & 0 \end{bmatrix}, \quad z = \begin{bmatrix} -1 & 1 \\ 1 & 1 \\ 1 & -1 \\ -1 & 0 \\ 1 & 0 \\ 1 & -1 \\ 1 & -1 \\ -1 & -1 \end{bmatrix}$$

and find a matrix  $\mathcal{H}$  that solves  $\mathbb{O}(x, z; \mathcal{T})$  without the constraint on  $\mathcal{M}_{x,z}$ . This will result in  $\mathcal{H}$  such that  $\mathcal{M}_{x,z}(\mathcal{H})$  has both negative and positive eigenvalues. For larger values of  $n$  and  $r$ , one can fill up the extra entries with zeros and the proof carries over.  $\square$

The code for reproducing the result is available on-line<sup>1</sup>. It took 482 tosses to generate the counterexample of the matrices containing only  $\pm 1$  or 0 as their components. We used the uniform distribution over those matrices to generate the tosses.

## 4.3 DC power systems with acyclic topology

In Section 4.2, we studied one particular structure for the operator  $\mathcal{A}$ . Now, we analyze a real-world problem to highlight the role of the KSP. Recall that the power

<sup>1</sup> [github.com/igormolybog/matrix-sense-global](https://github.com/igormolybog/matrix-sense-global)

system discussed in Section 2 was an AC network for which the voltages were complex numbers. To simplify the computation, we analyze a DC system in this section, where all voltages are real-valued (Ghosh, Boyd, and Saberi 2008). Assume that there are  $n$  nodes, associated with the unknown real-valued voltages  $\tilde{v}_1, \dots, \tilde{v}_n$ . The power is measured at each node  $i \in \{1, \dots, n\}$ , and is denoted as  $\tilde{p}_i$ , which can be calculated according to the formula:

$$p_i(\tilde{v}) = \sum_{j \in N(i)} \tilde{v}_i(\tilde{v}_i - \tilde{v}_j) \frac{1}{r_{ij}} = \tilde{p}_i,$$

where  $N(i) \subset \{1, \dots, n\}$  is the set of nodes adjacent to node  $i$  and  $r_{ij} = r_{ji} > 0$  is the resistance of the line between nodes  $i$  and  $j$ . The least-squares formulation of the voltage recovery problem consists in minimization over the set  $v \in \mathbb{B}_R(\mathbf{1})$  of

$$f(x) = \sum_{i=1}^n (p_i(v) - \tilde{p}_i)^2$$

which is a special case of the function (1). Let  $R$  be a number such that  $2v_i > v_n$  for all  $i \in \{1, \dots, n-1\}$ . In this subsection, we will demonstrate the application of our results on a specific topology of the network, although as discussed later on, our conclusion applies to any acyclic topology.

Suppose that the network possesses a star topology, meaning that each node  $i \in \{1, \dots, n-1\}$  is connected to only one node — namely, node  $n$  — and no others. This means that  $N(i) = \{n\}$  if  $i \neq n$  and  $N(n) = \{1, \dots, n-1\}$ . As a result, the power measurements in this particular case can be written as

$$\begin{aligned} p_i(v) &= v_i(v_i - v_n) \frac{1}{r_{in}}, \quad i \in \{1, \dots, n-1\} \\ p_n(v) &= \sum_{j=1}^{n-1} v_n(v_n - v_j) \frac{1}{r_{jn}} \end{aligned}$$

which generate a particular structure for the sensing operator. Solving  $a_i \text{vec}(vv^T - \tilde{v}\tilde{v}^T) = p_i$  for  $a_i$ , we conclude that the  $i$ -th row of  $\mathbf{A}$  can be written as

$$\begin{aligned} a_i &= \xi_i \text{vec}(E_{ii} - E_{ni}), \quad i \in \{1, \dots, n-1\} \\ a_n &= -\text{vec}\left(\xi_n E_{nn} + \sum_{j=1}^{n-1} \xi_j E_{jn}\right) \end{aligned}$$

where  $E_{ij}$  is an  $n \times n$  matrix with  $(i, j)$ -th entry equal to 1 and all other entries equal to zero, and where  $\xi_i = \frac{1}{r_{in}} = \frac{1}{r_{ni}} > 0$  for  $i \neq n$  and  $\xi_n = -\sum_{j=1}^{n-1} \xi_j$ . The corresponding kernel matrix is given by

$$\mathbf{H} = \mathbf{A}^T \mathbf{A} = \sum_{i=1}^n a_i a_i^T = \mathbf{H}_{\text{star}}(\xi)$$

As a result,  $\mathbf{H}$  has a structured sparsity pattern: it is a block-diagonal matrix with  $n$  blocks  $M_1, \dots, M_n \in \mathbb{S}^n$  such that the first  $n-1$  blocks have only four nonzero entries:

$$M_i = \xi_i^2 [E_{ii} - E_{in} - E_{ni} + E_{nn}] = \xi_i^2 [e_n - e_i][e_n - e_i]^T, \quad i \in \{1, \dots, n-1\}$$

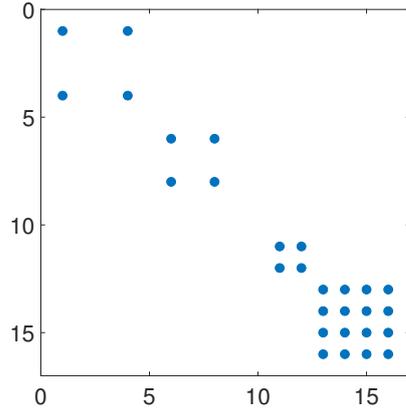


Fig. 2: Sparsity pattern of the matrix  $\mathbf{H}$  corresponding to a DC power system with a star topology consisting of four buses.

where  $e_i$  is the  $i$ -th column of the  $n \times n$  identity matrix. The last block of  $\mathbf{H}$  is a full matrix:

$$M_n = \xi \xi^T$$

The sparsity pattern of  $\mathbf{H}$  is visualized for  $n = 4$  in Figure 2. The matrix  $\mathbf{H} = \mathbf{H}_{star}(\xi)$  also has a low-dimensional structure: the  $4(n-1) + n^2$  non-zero entries of  $\mathbf{H}$  quadratically depend on only  $n-1$  parameters  $\xi_1, \dots, \xi_{n-1}$ . In Section 6, we will demonstrate how linearization of the structure can be applied, while here we provide an analytical proof that deals with the non-linear low-dimensional structure directly. This proof sheds light on some of the ideas behind Theorem 2.

The best RIP constant of the sensing operators that correspond to star topology power networks is significantly higher than  $\frac{1}{2}$ . The code repository mentioned earlier contains examples of  $v$  and  $\tilde{v}$  such that  $2v_i > v_n$  and  $2\tilde{v}_i > \tilde{v}_n$ , which prove that the best RIP constant of  $\mathcal{A}$  that correspond to a four-bus system with  $r_{14} = r_{24} = y_{34} = 1$  is at least 0.95. Therefore, Theorem 1 cannot be applied. Nevertheless, we will show that the non-convex voltage recovery problem on a system with a star topology possesses no spurious local minima.

**Proposition 3** Consider the problem

$$\left\{ \min_{x \in \mathbb{B}_{1/3}(\mathbf{1})} f_{z, \mathbf{A}}(x) \mid \mathbf{H} = \mathbf{H}_{star}(\xi); z \in \mathbb{B}_{1/3}(\mathbf{1}) \text{ and } \xi^T p(z) \neq 0 \right\},$$

or equivalently, (Problem<sup>KSP</sup>) under the additional constraint  $\xi^T p(z) \neq 0$ . No instance of this problem has a spurious second-order critical point.

*Proof.* Since  $R < \frac{1}{3}$ , it holds that  $2x_i > x_n$  and  $2z_i > z_n$  for all  $i \in \{1, \dots, n\}$ . We are interested in the landscape of the function

$$\begin{aligned} h(x) &= \text{vec}(xx^T - zz^T)^T \mathbf{H} \text{vec}(xx^T - zz^T) \\ &= \sum_{i=1}^n (x_i x - z_i z)^T M_i (x_i x - z_i z) \end{aligned}$$

To find the first and second derivatives, consider

$$\begin{aligned} h(x+u) &= \sum_{i=1}^n (x_i x - z_i z + x_i u + u_i x + u_i u)^T M_i (x_i x - z_i z + x_i u + u_i x + u_i u) \\ &= h(x) + 2 \sum_{i=1}^n (x_i u + u_i x)^T M_i (x_i x - z_i z) + \\ &\quad + \sum_{i=1}^n (x_i u + u_i x)^T M_i (x_i u + u_i x) + \\ &\quad + 2 \sum_{i=1}^n (u_i u)^T M_i (x_i x - z_i z) + o(|u|^2) \end{aligned}$$

Selecting the term that is linear in  $u$ , the gradient takes the form

$$\begin{aligned} \nabla_x h(x) &= 2 \sum_{i=1}^n [x_i M_i (x_i x - z_i z) + \text{Tr}[x^T M_i (x_i x - z_i z)] e_i] \\ &= 2 \sum_{i=1}^{n-1} \xi_i^2 x_i (e_n - e_i) (e_n - e_i)^T (x_i x - z_i z) + \\ &\quad + \text{Tr}[x^T (e_n - e_i) (e_n - e_i)^T (x_i x - z_i z)] e_i + \\ &\quad + 2x_i \xi_i \xi_i^T (x_i x - z_i z) + \text{Tr}[x^T \xi_i \xi_i^T (x_i x - z_i z)] e_i \end{aligned}$$

which can be written in the compact form

$$\nabla_x h(x) = B(x)[p(x) - p(z)]$$

where the  $(i, j)$ -th component of the  $n \times n$  matrix  $B(x)$  is

$$B_{ij} = \begin{cases} \xi_i (2x_i - x_n) & \text{if } i = j, i \neq n \\ \sum_{s=1}^{n-1} \xi_i (2x_n - x_s) & \text{if } i = j = n \\ -\xi_i x_i & \text{if } i \neq j, i = n \text{ or } i \neq j, j = n \\ 0 & \text{otherwise} \end{cases}$$

and  $p(x)$  is a vector with  $i$ -th component equal to  $p_i(x)$ . If  $B(x)$  is non-singular at a point  $x$ , then  $x$  is a first-order critical point of  $h(x)$  in the open ball  $\bar{\mathbb{B}}_{1/3}(\mathbf{1})$  if and only if  $p(x) - p(z) = \mathcal{A}(xx^T - zz^T) = 0$ , which implies that it is a global minimum. Therefore, it is essential to identify all points  $x$  such that  $\det(B) = 0$ .

With a slight abuse of notation, we show  $B(x)$  with the shorthand notation  $B$ . Let  $\beta_n$  denote the  $(n, n)$ -th entry of  $B$ , which is equal to  $\sum_{s=1}^{n-1} \xi_s (2x_n - x_s)$ . Represent the matrix  $B$  as a block matrix:  $B = \begin{bmatrix} B' & b^T \\ b & B'' \end{bmatrix}$  with the scalar  $B' = \xi_1 (2x_1 - x_n)$  and the  $(n-1)$ -dimensional vector  $b^T = [0, \dots, 0, -\xi_1 x_1]$ . Since  $x \in \bar{\mathbb{B}}_{1/3}(\mathbf{1})$ , we have  $B' \neq 0$ . One can write:

$$\det(B) = 0 \iff \det(B'' - B'^{-1} b b^T) = 0$$

The new  $(n-1) \times (n-1)$  matrix  $B'' - B'^{-1}bb^T$  is equal to  $B''$  in all components but its  $(n-1, n-1)$ -th entry, which changes to

$$\begin{aligned}\beta_{n-1} &= -\xi_1 \frac{x_1 x_n}{2x_1 - x_n} + \beta_n \\ &= -\xi_1 \frac{x_1 x_n}{2x_1 - x_n} + \xi_1 (2x_n - x_1) + \sum_{s=2}^{n-1} \xi_s (2x_n - x_s) \\ &= -2\xi_1 \frac{(x_1 - x_n)^2}{2x_1 - x_n} + \sum_{j=2}^{n-1} \xi_s (2x_n - x_s)\end{aligned}$$

Repeating the same matrix reduction argument  $n-1$  times yields that  $\det(B) = 0$  if and only if  $\beta_1 = 0$ , where

$$\beta_1 = -2 \sum_{i=1}^{n-1} \xi_i \frac{(x_i - x_n)^2}{2x_i - x_n}$$

Since  $2x_i - x_n > 0$ , we conclude that  $\beta_1 = 0$  if and only if  $x_1 = \dots = x_n$ . Therefore, there are no spurious solutions outside the set  $\{x : x_1 = \dots = x_n\}$ . To study this set, we derive the Hessian of  $h(x)$  by extracting from  $h(x+u)$  the term that is quadratic in  $u$ :

$$\begin{aligned}\nabla_{xx}^2 h(x) &= \sum_{i=1}^n [x_i^2 M_i + x_i (e_i x^T M_i + M_i x e_i^T) + e_i x^T M_i x e_i^T + \\ &\quad + M_i (x_i x - z_i z) e_i^T + e_i (x_i x - z_i z)^T M_i]\end{aligned}$$

and substitute  $x_1 = \dots = x_n = x'$  or  $x = x'\mathbf{1}$ . At the same time, we substitute  $M_i = \xi_i^2 [e_n - e_i][e_n - e_i]^T$  and  $M_n = \xi \xi^T$  and note that  $(e_n - e_i)^T \mathbf{1} = \mathbf{1}^T (e_n - e_i) = 0$  and  $\xi^T \mathbf{1} = \mathbf{1}^T \xi = 0$  by construction. After simplification, we obtain that

$$\begin{aligned}\nabla_{xx}^2 h(x)|_{x=x'\mathbf{1}} &= x'^2 \left[ \sum_{i=1}^{n-1} \xi_i^2 (e_n - e_i)(e_n - e_i)^T + \xi \xi^T - (\xi e_n^T + e_n \xi^T) z_n z^T \xi - \right. \\ &\quad \left. - \sum_{i=1}^{n-1} \xi_i^2 ((e_n - e_i) e_i^T + e_i (e_n - e_i)^T) z_i (z_n - z_i) \right]\end{aligned}$$

Consider the quadratic form  $q(s, t) = [s \ \dots \ s \ t] \nabla_{xx}^2 h(x)|_{x=x'\mathbf{1}} [s \ \dots \ s \ t]^T$ , where  $[s \ \dots \ s \ t] \in \mathbb{R}^n$ . One can write:

$$\begin{aligned}q(s, t) &= \sum_{i=1}^{n-1} \xi_i^2 (t-s)^2 + \left( \sum_{i=1}^{n-1} s \xi_i + t \xi_n \right)^2 - \\ &\quad - 2t z_n \left( \sum_{i=1}^{n-1} s \xi_i + t \xi_n \right) \left( \sum_{i=1}^{n-1} z_i \xi_i + z_n \xi_n \right) - \\ &\quad - 2s(t-s) \sum_{i=1}^{n-1} \xi_i^2 z_i (z_n - z_i) = \\ &= (t-s)^2 \left[ \sum_{i=1}^{n-1} \xi_i^2 + \left( \sum_{i=1}^{n-1} \xi_i \right)^2 \right] + \\ &\quad + 2(t-s) \left[ t \left( \sum_{i=1}^{n-1} \xi_i \right) \left( \sum_{i=1}^{n-1} \xi_i z_n (z_i - z_n) \right) - s \left( \sum_{i=1}^{n-1} \xi_i^2 z_i (z_n - z_i) \right) \right] = \\ &= c_1 (t-s)^2 + 2(t-s)(c_2 t - c_3 s) = \\ &= (c_1 + 2c_2)t^2 - 2(c_1 + c_2 + c_3)st + (c_1 + 2c_3)s^2\end{aligned}$$

where  $c_1, c_2$  and  $c_3$  are constants that are introduced to shorten the expression. If  $q(s, t)$  takes negative values, then  $\nabla_{xx}^2 h(x)|_{x=x^* \mathbf{1}}$  has a negative eigenvalue, and therefore  $x = x^* \mathbf{1}$  cannot be a spurious second-order critical point.

Now, consider the polynomials  $q(1, t)$  and  $q(s, 1)$ . Suppose that  $\sum_{i=1}^{n-1} \xi_i^2 \geq (\sum_{i=1}^{n-1} \xi_i)^2$  and consider any  $z \in \bar{\mathbb{B}}_{1/3}(\mathbf{1})$ . Since  $|z_n(z_i - z_n)| \leq \frac{4}{3} \cdot \frac{2}{3} = \frac{8}{9}$ , it must hold that

$$c_1 + 2c_2 > \sum_{i=1}^{n-1} \xi_i^2 + (\sum_{i=1}^{n-1} \xi_i)^2 - 2\frac{8}{9}(\sum_{i=1}^{n-1} \xi_i)^2 = \sum_{i=1}^{n-1} \xi_i^2 - \frac{6}{9}(\sum_{i=1}^{n-1} \xi_i)^2 \geq 0$$

Thus,  $q(1, t)$  is a polynomial of order 2 with respect to  $t$  with a positive leading term.

Suppose that  $\sum_{i=1}^{n-1} \xi_i^2 < (\sum_{i=1}^{n-1} \xi_i)^2$ . Similarly,

$$c_1 + 2c_3 > \sum_{i=1}^{n-1} \xi_i^2 + (\sum_{i=1}^{n-1} \xi_i)^2 - 2\frac{8}{9} \sum_{i=1}^{n-1} \xi_i^2 = (\sum_{i=1}^{n-1} \xi_i)^2 - \frac{6}{9} \sum_{i=1}^{n-1} \xi_i^2 \geq 0$$

and thus  $q(s, 1)$  is a polynomial of order 2 with respect to  $s$  with a positive leading term. At least one of the polynomials  $q(1, t)$  or  $q(s, 1)$  has a positive leading term. Both  $q(1, t)$  and  $q(s, 1)$  have the same determinant and thus the following argument can be made for any of them. Therefore, without loss of generality, assume that  $q(1, t)$  has a positive leading term and takes negative values if and only if its determinant is strictly positive:

$$4(c_1 + c_2 + c_3)^2 - 4(c_1 + 2c_2)(c_1 + 2c_3) > 0.$$

It is positive if and only if  $(c_2 - c_3)^2 > 0$ , which is equivalent to  $c_2 \neq c_3$ . After substituting  $c_2 = (\sum_{i=1}^{n-1} \xi_i)(\sum_{i=1}^{n-1} \xi_i z_n(z_i - z_n)) = -(\sum_{i=1}^{n-1} \xi_i) p_n(z)$  and  $c_3 = (\sum_{i=1}^{n-1} \xi_i^2 z_i(z_n - z_i)) = -\sum_{i=1}^{n-1} \xi_i p_i(z)$ , one can guarantee that  $\nabla_{xx}^2 h(x)|_{x=x^* \mathbf{1}}$  has a negative eigenvalue unless

$$\xi^T p(z) = 0.$$

Otherwise,  $q(1, t)$  only reaches zero at  $t = 1$  and never crosses it.  $\square$

The technique of using the properties of the Schur complement to eliminate the dimension of a matrix one by one can be applied in case of an arbitrary network with an acyclic topology. For any such network, the gradient  $\nabla_x h(x)$  also takes the form  $\nabla_x h(x) = B[p(x) - p(z)]$ , but with a different matrix  $B$ . Applying elimination of the rows and columns of  $B$  that correspond to the leaves first, then to the first layer of parent nodes, then to the second layer of parent nodes and so forth results in a similar result on the location of first-order critical points. Thus, in a similar way, the conclusion of Proposition 3 can be proven for any arbitrary acyclic network, but for a different value of the radius  $R$  that may not be analytically calculable. However, the proof is not generalizable to networks with cycles. The proof of Proposition 3 was based on analyzing the first- and second-order optimality conditions and exploiting the properties of the operator  $\mathcal{A}$  that benefits from both sparsity and a low-dimensional structure. The ideas used in the proof help the reader understand Theorem 2. Since Proposition 3 does not apply to networks with cycles, one may instead use Theorem 2 to numerically evaluate the inexistence of spurious solutions for any particular cyclic network. This will be carried out in Section 6.

## 5 Combining KSP with RIP

After fixing the hyperparameters  $\omega \in \mathbb{R}^{n \times r}$  and  $R \in \mathbb{R} \cup \{+\infty\}$  together with the kernel structure of the sensing operators and the RIP constant, we can state the problem under study in this section as follows:

$$\left\{ \min_{x \in \mathbb{B}_R(\omega)} f_{z, \mathcal{A}}(x) \mid \mathcal{A} \text{ satisfies } \delta_{2r}\text{-RIP and } \mathcal{T}\text{-KSP, } z \in \mathbb{B}_R(\omega) \right\},$$

(Problem<sup>KSP+RIP</sup>)

Note that (Problem<sup>KSP+RIP</sup>) consists in minimization of a class of functions  $f_{z, \mathcal{A}}$  that correspond to some point  $z \in \mathbb{B}_R(\omega)$  and some operator  $\mathcal{A}$  that satisfies  $\mathcal{T}$ -KSP and  $\delta_{2r}$ -RIP simultaneously. This is a generalization of both (Problem<sup>RIP</sup>) and (Problem<sup>KSP</sup>). For (Problem<sup>KSP+RIP</sup>), we provide necessary and sufficient conditions for having no spurious second-order critical point, and consequently no spurious local minimum.

**Theorem 3 (KSP+RIP necessary and sufficient conditions)** *For all instances of (Problem<sup>KSP+RIP</sup>), there are no spurious second-order critical points if*

$$\left\{ \delta_{2r} < \min_{\substack{x \in \mathbb{B}_R(\omega), z \in \mathbb{B}_R(\omega) \\ xx^T \neq zz^T}} \mathbb{O}_P(x, z; \mathcal{T}) \right. \quad (14a)$$

$$\left. \delta_{2r} < \min_{\substack{x \in \partial \mathbb{B}_R(\omega), z \in \mathbb{B}_R(\omega) \\ xx^T \neq zz^T}} \mathbb{O}_P^{\partial \mathbb{B}}(x, z; \mathcal{T}, \omega) \right. \quad (14b)$$

and only if

$$\left\{ \delta_{2r} < \min_{\substack{x \in \mathbb{B}_R(\omega), z \in \mathbb{B}_R(\omega) \\ xx^T \neq zz^T}} \mathbb{O}(x, z; \mathcal{T}) \right. \quad (15a)$$

$$\left. \delta_{2r} < \min_{\substack{x \in \partial \mathbb{B}_R(\omega), z \in \mathbb{B}_R(\omega) \\ xx^T \neq zz^T}} \mathbb{O}^{\partial \mathbb{B}}(x, z; \mathcal{T}, \omega) \right. \quad (15b)$$

Following from the results of Zhang, Sojoudi, and Lavaei 2019, the necessary and sufficient conditions coincide for the trivial structure operator  $\mathcal{T} \equiv 0$  and  $R = +\infty$ .

### 5.1 Robustness

Consider the scenario where the measurements are corrupted with independent and identically distributed Gaussian noise. More precisely, we assume that the measurement vector  $b$  is corrupted by an additive noise that can be written as  $\mathcal{A}(V)$  for some random matrix  $V$  that is probably full rank (since  $\mathcal{A}(\cdot)$  is from the high-dimensional space  $\mathbb{S}^n$  to the presumably low-dimensional space  $\mathbb{R}^m$ , we just need the mild surjectivity assumption). In this section, we show that the resulting recovery error can be bounded with high probability. For simplicity, we consider the case  $R = +\infty$ , but a similar argument can be used to analyse the case with a finite radius.

**Theorem 4** Consider (Problem<sup>KSP+RIP</sup>) with  $R = +\infty$  for which the condition (14a) holds. Let  $V \in \mathbb{S}^n$  be a random matrix of arbitrary rank. Define the noisy recovery loss

$$g(x) = \|\mathcal{A}(xx^T - zz^T + V)\|.$$

For every  $p \in (0, 1)$  and  $\varepsilon > 0$ , there exists  $\sigma = \sigma(p, \varepsilon; \mathcal{A}, z) > 0$  such that for  $V \sim \mathcal{N}(0, \sigma^2 I)$ , with probability at least  $p$ , every second-order critical point  $x^*$  of  $g(x)$  satisfies  $\|x^*x^{*T} - zz^T\| < \varepsilon$

*Proof.* Expand the recovery loss:

$$\begin{aligned} g(x) &= \langle xx^T - zz^T + V, (xx^T - zz^T + V) \rangle \\ &= f(x) + \langle V, \mathcal{H}(xx^T - zz^T) + xx^T - zz^T \rangle + \langle V, \mathcal{H}(V) \rangle \end{aligned}$$

We outline the proof below:

- Split  $\mathbb{R}^n$  into four regions according to the behaviour of  $f(x)$  associated with the noiseless scenario:
  1.  $\varepsilon$ -neighborhood of the second-order critical points of  $f(x)$
  2. some neighborhood of the remaining first-order critical points of  $f(x)$
  3. inner compact region where the value of  $\|\nabla f\|$  is bounded by some positive constants from below and from above
  4. Outer region, where  $\|\nabla f\|$  is large;
- Show the existence of  $\sigma$  such that there are no second-order critical points of  $g(x)$  in regions 2, 3 and 4 with high probability;
- Conclude that the only region that contains the second-order critical points of  $g(x)$  with high probability is region 1, which coincides with the set  $\{x : \|xx^T - zz^T\| < \varepsilon\}$ .

The illustration of the regions used in the proof can be found in Figure 3. To prove formally, first calculate the gradient

$$\nabla_x g(x) = \nabla_x f(x) + \nabla_x \langle V, \mathcal{H}(xx^T) + xx^T \rangle$$

For  $i \in \{1 \dots n\}$  and  $j \in \{1 \dots r\}$ , one can write:

$$\frac{\partial}{\partial x_{ij}} \langle V, xx^T \rangle = \langle V, e_i x_j^T + x_j e_i^T \rangle = 2e_i^T V x_j$$

where  $e_i$  is the  $i$ -th column of the  $n \times n$  identity matrix. Moreover,

$$\frac{\partial}{\partial x_{ij}} \langle V, \mathcal{H}(xx^T) \rangle = \langle V, \mathcal{H}(e_i x_j^T + x_j e_i^T) \rangle$$

The Hessian can also be written as

$$\nabla_{xx}^2 g(x) = \nabla_{xx}^2 f(x) + \nabla_{xx}^2 \langle V, \mathcal{H}(xx^T) + xx^T \rangle$$

Similarly, for  $i' \in \{1 \dots n\}$  and  $j' \in \{1 \dots r\}$ , we have

$$\frac{\partial^2}{\partial x_{ij} \partial x_{i'j'}} \langle V, xx^T \rangle = 2\delta_{jj'} \langle V, E_{ii'} \rangle$$

where  $\delta_{jj'} \in \mathbb{R}$  is defined as  $\delta_{jj'} = \begin{cases} 1 & \text{if } j = j' \\ 0 & \text{otherwise} \end{cases}$ , and  $E_{ii'} \in \mathbb{R}^{n \times n}$  is a matrix whose  $(i, i')$ -th entry is 1 and other entries are 0. Similarly,

$$\frac{\partial^2}{\partial x_{ij} \partial x_{i'j'}} \langle V, \mathcal{H}(xx^T) \rangle = \delta_{jj'} \langle V, \mathcal{H}(E_{ii'} + E_{i'i}) \rangle$$

By assumption, there exists  $\gamma > 0$  such that  $\frac{1-\delta_{2\epsilon}}{\gamma} \|xx^T - zz^T\|_F^2 \leq \|\mathcal{A}(xx^T - zz^T)\|_F^2$  for all  $x$ . This implies that  $f(x)$  is a coercive functions of  $x$  for any given  $\mathcal{A}$  and  $z$ . Moreover,  $\|\nabla_x f(x)\|$  is also a coercive function. To show this, using the notation from Section 4, consider

$$\begin{aligned} \langle \frac{x}{\|x\|_F}, \nabla_x f(x) \rangle &= \frac{2}{\|x\|_F} \langle \text{vec}(x), \mathbf{X}^T \mathbf{H} \mathbf{e} \rangle \\ &= \frac{2}{\|x\|_F} \langle \mathbf{X} \text{vec}(x), \mathbf{H} \mathbf{e} \rangle \\ &= \frac{2}{\|x\|_F} \langle xx^T + xx^T, \mathcal{H}(xx^T - zz^T) \rangle \\ &= \frac{4}{\|x\|_F} [f(x) + \langle zz^T, \mathcal{H}(xx^T - zz^T) \rangle] \end{aligned}$$

Knowing that  $f(x)$  grows as fast as  $\|xx^T\|_F^2 = \|x\|_F^4$  and  $-\langle zz^T, \mathcal{H}(xx^T - zz^T) \rangle$  grows at most as fast as  $\|x\|_F^2$ , we conclude that  $\langle \frac{x}{\|x\|_F}, \nabla_x f(x) \rangle \rightarrow \infty$  as  $x \rightarrow \infty$ , which implies that  $\|\nabla_x f(x)\| \rightarrow \infty$  as  $x \rightarrow \infty$ .

For an arbitrary  $K' > 0$ , define the set  $C_{K'} = \{x | f(x) \leq K', \|\nabla_x f(x)\| \leq K'\}$ . It is compact due to coerciveness. The difference  $\nabla_x g(x) - \nabla_x f(x)$  is linear in both  $V$  and  $x$ , while  $\nabla_x f(x)$  is cubic in  $x$ . Noting that  $\|\nabla_x g(x)\| \geq \|\nabla_x f(x)\| - \|\nabla_x g(x) - \nabla_x f(x)\|$ , one can conclude that  $\|\nabla_x g(x)\|$  is also a coercive function. Therefore, for any  $p_K \in (0, 1)$  there exist  $K$  and  $\sigma = \sigma_K$  such that  $\|\nabla_x g(x)\| > 0$  over  $\mathbb{R}^n \setminus C_K$  with probability  $p_K$ . Select  $p_K = \sqrt[3]{p}$  and fix the corresponding  $K$  and  $\sigma_K$ .

The set  $O_{fo}$  of first-order critical points of  $f(x)$  is closed due to the closed graph theorem. Moreover, it is bounded due to coerciveness of  $f(x)$ , and thus compact even when  $R = +\infty$ . Denote the set of second-order critical points of  $f(x)$  with  $O_{min} \subseteq O_{fo}$ . It coincides with the set of global minimizers of  $f(x)$  since the condition 14a holds and Theorem 3 can be utilized. Define  $U_{min} = \cup_{x \in O_{min}} \mathbb{B}_\epsilon(x)$ , and the set  $O_{rest} = O_{fo} \setminus U_{min}$  that is compact. Note that the minimal eigenvalue of  $\nabla_{xx}^2 f(x)$  is strictly negative for every  $x \in O_{rest}$ . Since minimal eigenvalue is a continuous function, there exists  $\bar{\lambda} < 0$  such that  $\min_{x \in O_{rest}} \lambda_{min}(\nabla_{xx}^2 f(x)) = \bar{\lambda}$ . By continuity of  $\nabla_{xx}^2 f(x)$  with respect to  $x$ , there exists  $\xi > 0$  such that  $\lambda_{min}(\nabla_{xx}^2 f(x)) < \frac{\bar{\lambda}}{2}$  for all  $x \in U_{rest} = \cup_{x' \in O_{rest}} \mathbb{B}_\xi(x')$ . The difference of Hessians  $\nabla_{xx}^2 g(x) - \nabla_{xx}^2 f(x)$  is linear in  $V$  and constant in  $x$ . Therefore, for any  $\psi > 0$  and  $p_\psi \in (0, 1)$ , there exists  $\sigma_\psi$  such that with probability  $p_\psi$ , it holds that  $\|\nabla_{xx}^2 g(x) - \nabla_{xx}^2 f(x)\|_F < \psi$  for all  $x \in C_K$ . Select  $\psi = \frac{\bar{\lambda}}{3}$  and  $p_\psi = \sqrt[3]{p}$ , and fix the corresponding  $\sigma_\psi$ . Notice that under  $\sigma \leq \sigma_\psi$ , with probability  $p_\psi$  there are no second-order critical points of  $g(x)$  in  $U_{rest}$ .

Denote  $U_{fo} = U_{rest} \cup U_{min}$  and notice that  $C_K \setminus U_{fo}$  is a compact set that contains no first-order critical points of  $f(x)$ . Therefore, there exists  $\rho > 0$  such that  $\|\nabla_x f(x)\| > \rho$  for all  $x \in C_K \setminus U_{fo}$ . Due to continuity of  $\nabla_x g(x) - \nabla_x f(x)$ , for any  $\phi > 0$  and  $p_\phi \in (0, 1)$ , there exists  $\sigma = \sigma_\phi$  such that with probability  $p_\phi$  it holds that  $\|\nabla_x g(x) - \nabla_x f(x)\|_F < \phi$  for all  $x \in C_K$ . Select  $\phi = \rho$  and  $p_\phi = \sqrt[3]{p}$ , and fix the corresponding

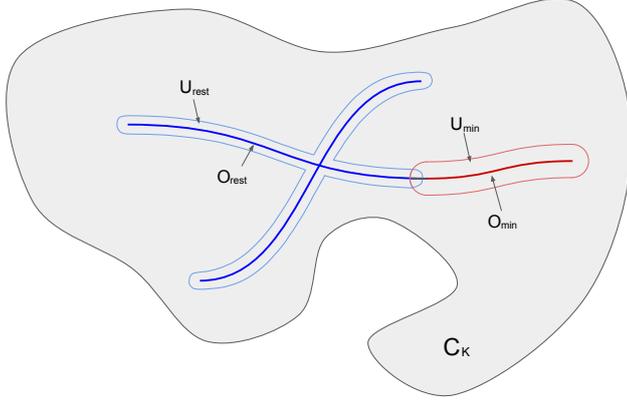


Fig. 3: Schematic of the domain of the function  $f(x)$  with highlighted regions. Grey area denotes the compact region  $C_K$ . The bold lines denote the set of first-order critical points named  $O_{fo}$  whose subset shown in red corresponds to the set of global minimizers named  $O_{min}$ , while the blue part corresponds to  $O_{rest}$ . The area countered by the red shaded line is the  $\varepsilon$ -neighborhood of  $O_{min}$ , namely  $U_{min}$ , while the area countered by the blue shaded line is the  $\xi$ -neighborhood of  $O_{rest}$ , namely  $U_{rest}$ . The proof finds that with high probability there are no second-order critical points of  $g(x)$  outside of  $C_K$  (outer region 4), or inside  $U_{rest}$  (region 2), or inside  $C_K \setminus [U_{rest} \cup U_{min}]$ . Therefore, all of such points must be located inside  $U_{min}$ .

$\sigma_\phi$ . Notice that under  $\sigma \leq \sigma_\phi$ , with probability  $p_\phi$  there are no second-order critical points of  $g(x)$  in  $C_K \setminus U_{fo}$ .

To conclude the proof, select  $\sigma < \min\{\sigma_K, \sigma_\psi, \sigma_\phi\}$  and observe that with probability at least  $p_K \times p_\psi \times p_\phi = p$  there are no second-order critical points of  $g(x)$  in the set  $\mathbb{R}^n \setminus U_{min} = [\mathbb{R}^n \setminus C_K] \cup U_{rest} \cup [C_K \setminus U_{fo}]$ .

□

## 5.2 Sparse structure and normalization

Due to Theorem 1 for the rank-1 case, the instances of (Problem<sup>KSP+RIP</sup>) have no spurious solutions with  $\mathcal{T} \equiv 0$  as long as  $\delta_2$  is upper bounded by  $\frac{1}{2}$ . In this subsection, we are concerned with the question of how much sparsity can impact the best bound on RIP that certifies global convergence. Formally, we set  $\mathcal{W} \equiv 0$  and  $\mathcal{T} \equiv \mathcal{S}$  and find a tighter upper bound for  $\delta_2$ . After enforcing sparsity, it is natural to expect that the bound grows and becomes less restrictive. However, this turns out not to be the case.

Let  $n = 2$  and  $r = 1$ , and consider the smallest sparsity pattern possible for  $\mathcal{H} = \mathcal{A}^T \mathcal{A} \succ 0$ . It consists exclusively of elements  $(i, i)$ , and thus enforces  $\mathbf{H}$  to be diagonal. Consider the point  $x$  with respect to the instance of the problem given by  $z$  and  $\mathbf{A}$  as in the example below:

*Example 1* Assume that

$$x = (1, 1); \quad z = (\sqrt{2}, -\sqrt{2}); \quad \mathbf{A} = \text{diag}(\sqrt{3}, 1, 1, \sqrt{3})$$

Then,  $x$  is spurious for  $f_{z, \mathbf{A}}$  since it satisfies the second-order necessary conditions:

$$\nabla f_{z, \mathbf{A}}(x) = 0, \quad \nabla^2 f_{z, \mathbf{A}}(x) = 16 \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \succeq 0$$

which makes it a spurious second-order critical point (note that  $xx^T \neq zz^T$ ). Notice that  $\mathcal{H} = \mathcal{A}^T \mathcal{A}$  is indeed diagonal. Moreover, for all  $X \in \mathbb{S}^2$ , the operator  $\mathcal{A}$  satisfies the tight bound  $\|X\|_F^2 \leq \|\mathcal{A}(X)\|^2 = \left\| \begin{bmatrix} \sqrt{3} & 1 \\ 1 & \sqrt{3} \end{bmatrix} \circ X \right\| \leq 3\|X\|_F^2$ . Therefore, the largest number  $\delta_2$  for this instance is equal to  $1/2$ , which coincides with the upper bound for unstructured problems. Somewhat counter-intuitively, the tight bound established in Zhang et al. 2018; Zhang, Sojoudi, and Lavaei 2019 holds even when a very restrictive sparsity pattern of the kernel operator is enforced. Nevertheless, for an arbitrary low-dimensional structure  $\mathcal{W}$ , a tighter sparsity constraint entails a less restrictive bound on incoherence as discussed below.

**Proposition 4** *If the sparsity pattern  $S$  has a sub-pattern  $S'$  meaning that  $S' \subset S$ , then  $\mathbb{O}(x, z; \mathcal{W}, \mathcal{S}') \leq \mathbb{O}(x, z; \mathcal{W}, \mathcal{S})$  for all  $x, y \in \mathbb{R}^{n \times r}$ . Thus, the necessary bound on incoherence for  $\mathbf{H}$  with  $S'$  is not more restrictive than the bound for  $\mathbf{H}$  with  $S$ .*

In other words, a more restricting assumption on the sparsity of the kernel operator can only push the upper bound on the RIP constant higher up. Consequently, Example 1 shows that there is no sparsity pattern of cardinality  $> 3$  that can itself compensate the lack of isometry. Note that the example is given for the case  $n = 2$ , but there is a straightforward extension to an arbitrary  $n$  by adding zero components to  $x$  and  $z$ . It is common in practice to normalize the rows of the sensing matrix before proceeding to recovery. In the context of power systems, it is expressed as  $x^T M_i x \rightarrow \frac{x^T M_i x}{\|M_i\|_F}$ . For Example 1, after normalization,  $\mathbf{A}$  turns into the identity. The corresponding instance of the problem is known to have no spurious critical points. This illustrates how normalization helps to improve the isometry property of the sensing operator and removes the spurious second-order critical points out of the corresponding instance of the problem. Normalization in this case can be regarded as inducing structure on top of sparsity.

## 6 Numerical results

It is desirable to numerical study the non-convex matrix recovery in problems with a structured sensing operator. The objective is to show how the general theory developed in Section 3 can be applied to a real-world problem, namely the power system

state estimation discussed in Section 2. In general, optimization problems in (14) and (15) are non-convex. Thus, we propose to use Bayesian optimization (Frazier 2018) in order to obtain a numerical estimation of their solutions.

### 6.1 Power systems

In this section, we focus our attention on three networks named `case9`, `case14` and `case30` that are provided in the MATPOWER package. For `case9`, the number of buses is  $n = 9$  and there are  $m = 63$  possible power measurements that can be collected, while we have  $n = 14$  and  $m = 98$  for `case14` and  $n = 30$  and  $m = 210$  for `case30`. We denote the corresponding sensing operators with  $\mathcal{A}^9$ ,  $\mathcal{A}^{14}$  and  $\mathcal{A}^{30}$ . Although the matrices  $\mathbf{A}$  have complex entries, the corresponding kernel operators  $\mathcal{H}$  are represented with real matrices due to the properties of the measurement matrices. Both matrices  $\mathbf{A}^{30}$  and  $\mathbf{H}^{30}$  are visualized in Figure 1.

We linearize the low-dimensional structure that was discussed in Section 4.3. Repetition of the non-zero entries of  $\mathbf{H}$  (after some scaling) is considered as a form of low-dimensional structure, instead of the nonlinear dependence on the admittance. For example, if the entries  $(i, j)$  and  $(i', j')$  of  $\mathcal{H}^{30}$  are equal, then  $\mathcal{W}^{30}$  is constructed to be such that its kernel consists of matrices, for which the entries  $(i, j)$  and  $(i', j')$  are equal.

Based on this property, we form the linear operators  $\mathcal{T}^9$ ,  $\mathcal{T}^{14}$  and  $\mathcal{T}^{30}$ . All of the matrices in their kernel subspace are rank deficient. In this case, Theorem 3 can only provide us with the trivial upper bound on the RIP:  $\delta_2 < 1$ . However, this operator will allow us to use Theorem 3 to find a less conservative bound on RIP to certify the inexistence of spurious solutions for the structured mapping. Although the power system state estimation aims to find a complex vector, it is straightforward to verify that  $\langle a + \sqrt{-1}b, \mathbf{H}(a + \sqrt{-1}b) \rangle = \langle a, \mathbf{H}a \rangle + \langle b, -\mathbf{H}b \rangle$  for any real vectors  $a$  and  $b$  as well as a real symmetric matrix  $\mathbf{H}$ . Therefore, it is enough to consider (14) over  $\mathcal{X} = \mathcal{Z} = \mathbb{R}^n$ .

Purpose of the experiment is to study the dependence of  $\delta_{2r}$  that is sufficient for absence of spurious solutions in (Problem<sup>KSP+RIP</sup>) on the radius  $R$  of the ball domain. Intuitively, one would expect the dependence to be monotonically decreasing, since the larger the domain, the more solutions can appear there with some being spurious. However, this is not exactly what can be observed. Figure 4 shows the right-hand side of the inequalities in (14) from Theorem 3 for a range of values of  $R$  for three structure operators:  $\mathcal{T}^9$ ,  $\mathcal{T}^{14}$  and  $\mathcal{T}^{30}$ . In these experiments, the vector  $\omega$  has the unit entries. The red line provides with a guarantee on no spurious solution in the interior of the domain, while the blue dashed line takes care of the spurious solutions on the boundary. Indeed, the red curve decreases monotonically and converges to a value around 0.64 for all the experiments, while the blue dashed line decreases to 0.5 and recovers back to the same value afterwards. It turns out that 0.64 is the bound on  $\delta_{2r}$  for  $R = +\infty$  in each of the cases as well. This interesting behaviour is possible to be explained qualitatively. Consider a toy example with three cases in Figure 5, where the domain grows from Case I to Case III. There are no spurious solutions in case I, whereas one appears in case II and disappears in case III. Notice that the spurious

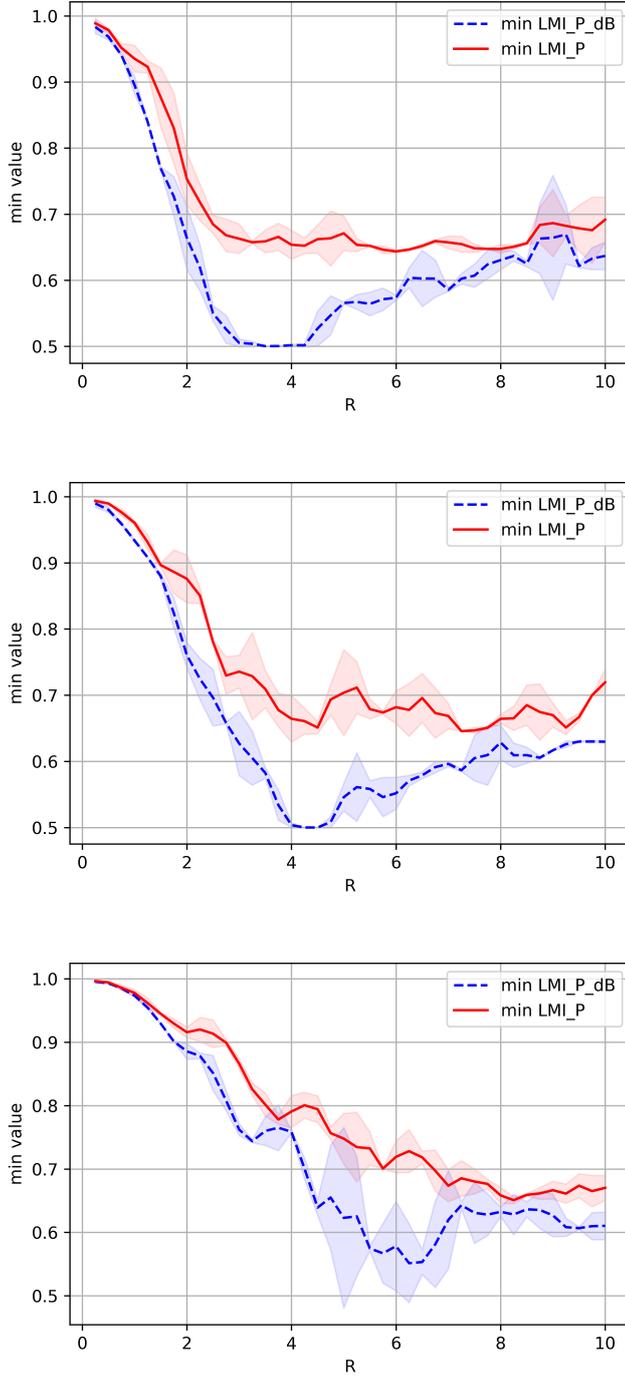


Fig. 4: The outcome of the minimization of  $\mathbb{O}_P(x, z)$  and  $\mathbb{O}_P^{\partial \mathbb{B}}(x, z)$  with the Bayesian optimization toolbox. The resulting value is the approximation of the right-hand side of the inequalities in (14) and can be used in Theorem 3 to estimate the lower bound on the sufficient RIP constant for global optimality. The values of the radius of the domain ball  $\mathbb{B}_R(\omega)$  are on the x-axis, and the corresponding approximations of  $\min \mathbb{O}_P(x, z)$  and  $\min \mathbb{O}_P^{\partial \mathbb{B}}(x, z)$  are on the y-axis. The red line depicts the lowest observed value of the function  $\mathbb{O}_P(x, z)$  and the blue dashed line depicts the minimum value of the function  $\mathbb{O}_P^{\partial \mathbb{B}}(x, z)$ .

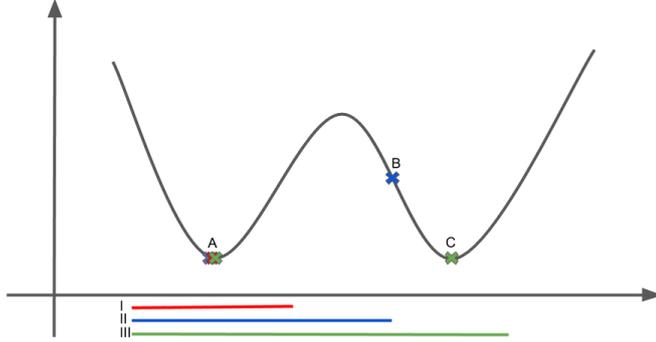


Fig. 5: Illustration for the local solution on the boundary. Three cases are considered, each marked with a different color. The colored intervals along the  $x$ -axis depict the domain in each of the cases, while the colored crosses denote the local solutions.

solution can only appear on the boundary, which motivates the steady behaviour of the red curve in Figure 4. Recall that the threshold 0.5 is valid for the trivial structure operator  $\mathcal{T} \equiv 0$  and  $R = +\infty$  and the blue curve never goes below it. Therefore, the constructed conditions of absence of spurious local optimality are strictly superior to the previously known bound.

The above simulations were based on the networks provided in the package MATPOWER 7.0b1 (Zimmerman, Murillo-Sánchez, and Thomas 2011). Keeping the structure of a network, we set the parameters of the lines equal to each other to be able to better visualize the operator  $\mathcal{H}$ . All the presented simulations were done using the MATLAB bayesopt toolbox, and MATLAB modeling toolbox CVX (Grant and Boyd 2014, 2008) with SDPT3 (Toh, Todd, and Tütüncü 1999; Tütüncü, Toh, and Todd 2003) as the underlying solver.

## 6.2 Synthetic data

In this subsection, we present numerical studies of the matrix recovery problem for structured sensing operators obtained from random ensembles. For simplicity, we set  $R = +\infty$ . Here, we propose to use Bayesian optimization (Frazier 2018) in order to obtain a numerical estimation of the solution of the optimization problem (14). We have empirically observed that Bayesian optimization tends to obtain the same optimal solution to this problem much faster than random shooting or cross-entropy. In this section, the smallest value of  $\delta_r$  such that  $\mathcal{A}$  satisfies the  $\delta_r$ -RIP property is referred to as *the best RIP constant* of the map  $\mathcal{A}$ .

Recall that the structure operator is defined by two operators stack together:  $\mathcal{T} = (\mathcal{S}, \mathcal{W})$ . Here,  $\mathcal{W}$  captures the underlying structure that is not captured by the sparsity operator  $\mathcal{S}$ . We consider the same form of this operator as in the experiment on power systems data. Given the matrix representation  $\mathbf{H}$  of the kernel operator, de-

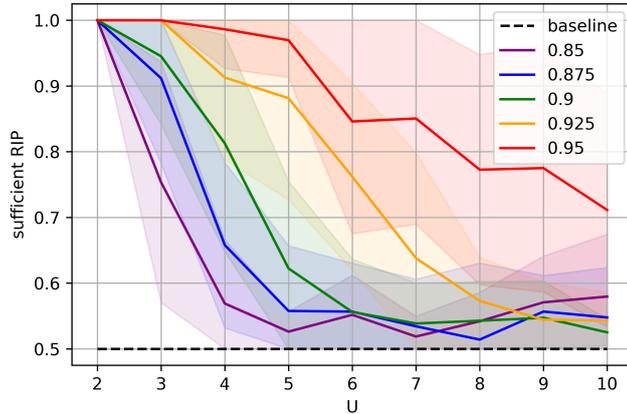


Fig. 6: The average of sufficient best RIP constant obtained from the developed analytic framework (Theorem 3) for random structures generated from the distribution  $RS(p_0, U)$  (each colored line stands for one specific value of  $p_0$ ), compared with the baseline method from Theorem 1 (shown as black and dashed). Shaded area represents the standard deviation window.

note the unique nonzero values in this matrix with the scalars  $h_1, \dots, h_{d_{\mathcal{W}}}$ . It means that  $\mathbf{H}$  is representable in the form  $\mathbf{H} = h_1 \mathbf{E}_1 + \dots + h_{d_{\mathcal{W}}} \mathbf{E}_{d_{\mathcal{W}}}$ , where  $\mathbf{E}_i$  is a matrix of the same size as  $\mathbf{H}$ , with 0 and 1 entries. The operator  $\mathcal{W}$  that we use in this section is any operator that has the subspace  $\{\beta_1 \mathbf{E}_1 + \dots + \beta_{d_{\mathcal{W}}} \mathbf{E}_{d_{\mathcal{W}}} \mid \beta_1, \dots, \beta_{d_{\mathcal{W}}} \in \mathbb{R}\}$  as its kernel.

We introduce a distribution  $RS(p_0, U)$  over the space of structure operators by describing the sampling scheme below. First, we generate the measurement structure matrix  $\mathbf{A}_{\text{st}}$  such that each of its components takes value 0 with probability  $p_0$  and any of the values  $1, \dots, U$  with the equal probability of  $\frac{1-p_0}{U}$ . We then form the kernel structure matrix as  $\mathbf{H}_{\text{st}} = \mathbf{A}_{\text{st}}^T \mathbf{A}_{\text{st}}$  and construct the sparsity operator  $\mathcal{S}$  and the extra structure operator  $\mathcal{W}$  as discussed before. The obtained structure operator  $\mathcal{T}$  is such that the operator represented with  $\mathbf{A}_{\text{st}}$  satisfies the  $\mathcal{T}$ -KSP. Note that the average sparsity of  $\mathbf{A}_{\text{st}}$  is  $p_0$  and the number of unique nonzero values is  $U$  with high probability, which implies that  $p_0$  is the parameter for the amount of sparsity structure in the problem, and  $U$  is the parameter for the amount of additional structure.

Figure 6 depicts the estimated sufficient RIP to guarantee the existence of no spurious second-order critical points, random problems with different values for the sparsity ( $p_0$ ) and the unique counter ( $U$ ). The sufficient RIP is obtained from Theorem 3 by imposing the KSP. Observe that the sparsity and the additional structure (the number of unique nonzero values in the measurement matrix in this particular case) both have a significant impact on the sufficient RIP. Note that higher  $p_0$  means more sparsity and lower  $U$  means more extra structure. Although it was observed theoretically that sparsity alone cannot guarantee the increase in the sufficient best RIP

constant, it appears to be an important characteristic when combined with the additional structure. Even for structures with a considerably low sparsity (0.85), the tight extra structure ( $U = 2$ ) has the sufficient best RIP of 1, which is a counter-intuitive result. The sufficient RIP seems to decay exponentially as we relax extra structure by increasing  $U$ , but with different bases for different  $p_0$ . This behaviour coincides with the one predicted in Proposition 4. If the goal is to make the RIP higher than a certain threshold, the amount of extra structure needed to achieve this reduces dramatically with the increase of the sparsity structure.

The experiment demonstrates that our method can be successfully applied to matrix sensing with randomly generated structure. The key takeaway from this experiment is that our method captures the trade-off between the sparsity and the low-dimensional structural properties of a given mapping. It shows that imposing restrictions on structure significantly affects the sufficient RIP, which leads to certifying the absence of spurious solutions under far less restrictive requirements (by improving the previous RIP bound 0.5 for arbitrary mappings).

## 7 Conclusion

The paper is concerned with the theoretical explanation of the recent empirical success of solving the low-rank matrix sensing problem via nonconvex optimization. It is known that under an incoherence assumption (namely, RIP) on the sensing operator, the optimization problem has no spurious local minima. This assumption is too strong for real-world applications where the amount of data cannot be sufficiently high. Aside from that, it does not account for the prior about the solution that is available in different applications. We develop the notion of Kernel Structure Property (KSP) based on linear matrix inequalities, which can be used instead or combined with RIP in this context. KSP explains how the inherent structure of an operator contributes to the inexistence of spurious local minima over the entire space or a given ball. As a special case, we study sparse sensing operators that have a low-dimensional representation. Using KSP, we obtain novel necessary and sufficient conditions for having no spurious solutions over a compact set for the matrix sensing problem, and demonstrate them in analytical and numerical studies.

## References

- Agarwal, Alekh, et al. 2016. "Learning sparsely used overcomplete dictionaries via alternating minimization". *SIAM Journal on Optimization* 26 (4): 2775–2799.
- Bhojanapalli, Srinadh, Behnam Neyshabur, and Nati Srebro. 2016. "Global optimality of local search for low rank matrix recovery". In *Advances in Neural Information Processing Systems*, 3873–3881.
- Bottou, Léon, and Olivier Bousquet. 2008. "The tradeoffs of large scale learning". In *Advances in neural information processing systems*, 161–168.
- Candes, Emmanuel J, Xiaodong Li, and Mahdi Soltanolkotabi. 2015. "Phase retrieval via Wirtinger flow: Theory and algorithms". *IEEE Transactions on Information Theory* 61 (4): 1985–2007.

- Chen, Yuxin, et al. 2018. “Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval”. *Mathematical Programming*: 1–33.
- Chen, Yuxin, et al. 2019. “Noisy matrix completion: understanding statistical guarantees for convex relaxation via nonconvex optimization”. *arXiv preprint arXiv:1902.07698*.
- Chi, Yuejie, Yue M Lu, and Yuxin Chen. 2018. “Nonconvex optimization meets low-rank matrix factorization: An overview”. *arXiv preprint arXiv:1809.09573*.
- Daneshmand, Hadi, et al. 2018. “Escaping saddles with stochastic gradients”. *arXiv preprint arXiv:1803.05999*.
- Frazier, Peter I. 2018. “A tutorial on bayesian optimization”. *arXiv preprint arXiv:1807.02811*.
- Ge, Rong, Chi Jin, and Yi Zheng. 2017. “No spurious local minima in nonconvex low rank problems: A unified geometric analysis”. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, 1233–1242. JMLR.org.
- Ge, Rong, Jason D Lee, and Tengyu Ma. 2016. “Matrix completion has no spurious local minimum”. In *Advances in Neural Information Processing Systems*, 2973–2981.
- Ge, Rong, et al. 2015. *Escaping From Saddle Points — Online Stochastic Gradient for Tensor Decomposition*. arXiv: 1503.02101 [cs.LG].
- Ghosh, Arpita, Stephen Boyd, and Amin Saberi. 2008. “Minimizing effective resistance of a graph”. *SIAM review* 50 (1): 37–66.
- Grant, Michael, and Stephen Boyd. 2014. *CVX: Matlab Software for Disciplined Convex Programming, version 2.1*.
- . 2008. “Graph implementations for nonsmooth convex programs”. In *Recent Advances in Learning and Control*, ed. by V. Blondel, S. Boyd, and H. Kimura, 95–110. Lecture Notes in Control and Information Sciences. Springer-Verlag Limited.
- Jain, Prateek, Praneeth Netrapalli, and Sujay Sanghavi. 2013. “Low-rank matrix completion using alternating minimization”. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, 665–674. ACM.
- Josz, Cedric, et al. 2018. “A theory on the absence of spurious solutions for nonconvex and nonsmooth optimization”. In *Advances in neural information processing systems*, 2441–2449.
- Keshavan, Raghunandan H, Andrea Montanari, and Sewoong Oh. 2010a. “Matrix completion from a few entries”. *IEEE transactions on information theory* 56 (6): 2980–2998.
- . 2010b. “Matrix completion from noisy entries”. *Journal of Machine Learning Research* 11 (Jul): 2057–2078.
- Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E Hinton. 2012. “Imagenet classification with deep convolutional neural networks”. In *Advances in neural information processing systems*, 1097–1105.
- Lee, Jason D, et al. 2016. “Gradient descent only converges to minimizers”. In *Conference on learning theory*, 1246–1257.
- Li, Dawei, Tian Ding, and Ruoyu Sun. 2018. “Over-Parameterized Deep Neural Networks Have No Strict Local Minima For Any Continuous Activations”. *arXiv preprint arXiv:1812.11039*.

- Li, Qiuwei, Zhihui Zhu, and Gongguo Tang. 2019. “Alternating minimizations converge to second-order optimal solutions”. In *International Conference on Machine Learning*, 3935–3943.
- Li, Xingguo, et al. 2019. “Symmetry, saddle points, and global optimization landscape of nonconvex matrix factorization”. *IEEE Transactions on Information Theory*.
- Mei, Song, Yu Bai, and Andrea Montanari. 2016. “The landscape of empirical risk for non-convex losses”. *arXiv preprint arXiv:1607.06534*.
- Nouiehed, Maher, and Meisam Razaviyayn. 2018. “Learning deep models: Critical points and local openness”. *arXiv preprint arXiv:1803.02968*.
- Pardalos, Panos M., and Stephen A. Vavasis. 1991. “Quadratic programming with one negative eigenvalue is NP-hard”. *Journal of Global Optimization* 1, no. 1 (): 15–22. ISSN: 1573-2916. doi:10.1007/BF00120662.
- Park, Dohyung, et al. 2016. “Non-square matrix sensing without spurious local minima via the Burer-Monteiro approach”. *arXiv preprint arXiv:1609.03240*.
- Paternain, Santiago, Aryan Mokhtari, and Alejandro Ribeiro. 2019. “A newton-based method for nonconvex optimization with fast evasion of saddle points”. *SIAM Journal on Optimization* 29 (1): 343–368.
- Soltanolkotabi, Mahdi. 2017. “Learning relus via gradient descent”. In *Advances in Neural Information Processing Systems, 2007–2017*.
- Sun, Ju, Qing Qu, and John Wright. 2018. “A geometric analysis of phase retrieval”. *Foundations of Computational Mathematics* 18 (5): 1131–1198.
- . 2015. “Complete dictionary recovery using nonconvex optimization”. In *International Conference on Machine Learning*, 2351–2360.
- Sun, Ruoyu, and Zhi-Quan Luo. 2016. “Guaranteed matrix completion via non-convex factorization”. *IEEE Transactions on Information Theory* 62 (11): 6535–6579.
- Toh, Kim-Chuan, Michael J Todd, and Reha H Tütüncü. 1999. “SDPT3: MATLAB software package for semidefinite programming, version 1.3”. *Optimization methods and software* 11 (1-4): 545–581.
- Tütüncü, Reha H, Kim-Chuan Toh, and Michael J Todd. 2003. “Solving semidefinite-quadratic-linear programs using SDPT3”. *Mathematical programming* 95 (2): 189–217.
- Vaswani, Namrata, Seyedehsara Nayer, and Yonina C Eldar. 2017. “Low-rank phase retrieval”. *IEEE Transactions on Signal Processing* 65 (15): 4059–4074.
- Yun, Chulhee, Suvrit Sra, and Ali Jadbabaie. 2018. “Global Optimality Conditions for Deep Neural Networks”. In *International Conference on Learning Representations*.
- Zhang, Richard Y., Somayeh Sojoudi, and Javad Lavaei. 2019. “Sharp Restricted Isometry Bounds for the Inexistence of Spurious Local Minima in Nonconvex Matrix Recovery”. *Journal of Machine Learning Research* 20 (114): 1–34.
- Zhang, Richard, et al. 2018. “How much restricted isometry is needed in nonconvex matrix recovery?” In *Advances in neural information processing systems*, 5591–5602.
- Zhang, Yu, Ramtin Madani, and Javad Lavaei. 2018. “Conic relaxations for power system state estimation with line measurements”. *IEEE Transactions on Control of Network Systems* 5 (3): 1193–1205.

- Zhao, Tuo, Zhaoran Wang, and Han Liu. 2015. “A nonconvex optimization framework for low rank matrix estimation”. In *Advances in Neural Information Processing Systems*, 559–567.
- Zheng, Qinqing, and John Lafferty. 2015. “A convergent gradient descent algorithm for rank minimization and semidefinite programming from random linear measurements”. In *Advances in Neural Information Processing Systems*, 109–117.
- Zhu, Zhihui, et al. 2018. “Global optimality in low-rank matrix optimization”. *IEEE Transactions on Signal Processing* 66 (13): 3614–3628.
- Zimmerman, Ray Daniel, Carlos Edmundo Murillo-Sánchez, and Robert John Thomas. 2011. “MATPOWER: Steady-state operations, planning, and analysis tools for power systems research and education”. *IEEE Transactions on power systems* 26 (1): 12–19.

## Appendix

In this part, we will prove Theorems 2 and 3 by nonexistence of a counterexample. Specifically, given  $\mathcal{T}$ , for a point  $x$  and a parameter value  $z$ , we aim to find a value  $\delta_{2r}^{x,z}$  for which the following claim holds:

“There exists  $\mathcal{A}$  that satisfies  $\mathcal{T}$ -KSP and  $\delta_{2r}$ -RIP such that  $x$  is a second-order critical point of  $f_{z,\mathcal{A}}$  if and only if  $\delta_{2r} > \delta_{2r}^{x,z}$ ”

$\mathcal{X}$  is equal to  $\mathbb{B}_R(\omega)$  in the notation from Section 1. The conditions for a point  $x$  to be a second-order critical point of a function  $f$  over  $\mathbb{B}_R(\omega)$  can be expressed in the compact form:

$$\begin{cases} \nabla f(x) = 0, \\ \nabla^2 f(x) \succeq 0 \end{cases} \quad \text{if } x \notin \partial \mathbb{B}_R(\omega) \text{ or } \begin{cases} \exists \mu \leq 0 : \nabla f(x) = \mu(x - \omega), \\ P_{x-\omega} \nabla^2 f(x) P_{x-\omega}^T \succeq 0 \end{cases} \quad \text{if } x \in \partial \mathbb{B}_R(\omega)$$

where  $P_{x-\omega} \in \mathbb{R}^{(nr-1) \times nr}$  is the matrix of orthogonal projection onto the subspace orthogonal to  $x - \omega$ . With that in mind, we construct two functions:  $\delta(x, z)$  and  $\partial \delta(x, z)$  by the following optimization procedures:

$$\begin{aligned} \delta(x, z) &\equiv \underset{\delta_{2r} \in \mathbb{R}, \mathcal{A}}{\text{minimum}} \quad \delta_{2r} \\ &\text{subject to} \quad \mathcal{L}_{x,z}(\mathcal{A}^T \mathcal{A}) = 0 \\ &\quad \mathcal{M}_{x,z}(\mathcal{A}^T \mathcal{A}) \succeq 0 \\ &\quad \mathcal{T}(\mathcal{A}^T \mathcal{A}) = 0 \\ &\quad \mathcal{A} \text{ satisfies } \delta_{2r}\text{-RIP.} \end{aligned}$$

$$\begin{aligned} \partial \delta(x, z) &\equiv \underset{\delta_{2r}, \mu \in \mathbb{R}, \mu \geq 0, \mathcal{A}}{\text{minimum}} \quad \delta_{2r} \\ &\text{subject to} \quad \mathcal{L}_{x,z}(\mathcal{A}^T \mathcal{A}) = -\mu(x - \omega) \\ &\quad P_{x-\omega} \mathcal{M}_{x,z}(\mathcal{A}^T \mathcal{A}) P_{x-\omega}^T \succeq 0 \\ &\quad \mathcal{T}(\mathcal{A}^T \mathcal{A}) = 0 \\ &\quad \mathcal{A} \text{ satisfies } \delta_{2r}\text{-RIP.} \end{aligned}$$

In each of the problems, the first two constraints represent the requirement that  $x$  is a second-order critical point is  $f_{z,\mathcal{A}}$ , the third constraint takes care of the KS property, and the last one of the RIP. It is straightforward to verify that  $\min\{\delta, \partial\delta\}$  takes the value of the desired  $\delta_{2r}^{x,z}$ . Minimization of  $\delta_{2r}^{x,z}$  over  $\{x \in \mathcal{X}, z \in \mathcal{Z} : xx^T \neq zz^T\}$  gives  $\delta_{2r}^*$  such that (Problem<sup>KSP+RIP</sup>) with  $\delta_{2r}$  has an instance with a spurious second-order critical point if and only if  $\delta_{2r} > \delta_{2r}^*$ .

Suppose that we are able to find  $\underline{\delta}_{2r}^{x,z}$  and  $\overline{\delta}_{2r}^{x,z}$  such that  $\underline{\delta}_{2r}^{x,z} \leq \delta_{2r}^{x,z} \leq \overline{\delta}_{2r}^{x,z}$  for all  $x \in \mathcal{X}, z \in \mathcal{Z}$ . Then,

$$\delta^* = \min_{\substack{x \in \mathcal{X}, z \in \mathcal{Z} \\ xx^T \neq zz^T}} \delta_{2r}^{x,z} \leq \min_{\substack{x \in \mathcal{X}, z \in \mathcal{Z} \\ xx^T \neq zz^T}} \underline{\delta}_{2r}^{x,z} \leq \min_{\substack{x \in \mathcal{X}, z \in \mathcal{Z} \\ xx^T \neq zz^T}} \overline{\delta}_{2r}^{x,z} = \overline{\delta}^*.$$

This inequality shows that  $\delta_{2r} \geq \underline{\delta}_{2r}^*$  is a sufficient, and  $\delta_{2r} \leq \overline{\delta}_{2r}^*$  is a necessary condition for the absence of spurious second-order critical points in the instances of the problem (Problem<sup>KSP+RIP</sup>). Now, it is desirable to show that  $\min\{\mathbb{O}_P^{\partial\mathbb{B}}(x, z; \mathcal{T}, \omega), \mathbb{O}_P(x, z; \mathcal{T})\}$  can serve as  $\underline{\delta}_{2r}^{x,z}$ , and  $\min\{\mathbb{O}^{\partial\mathbb{B}}(x, z; \mathcal{T}, \omega), \mathbb{O}(x, z; \mathcal{T})\}$  as  $\overline{\delta}_{2r}^{x,z}$ .

**Lemma 4** *The following statements hold all  $x \in \mathcal{X}$  and  $z \in \mathcal{Z}$ :*

$$\mathbb{O}_P(x, z) \leq \delta(x, z) \leq \mathbb{O}(x, z) \quad (16a)$$

$$\mathbb{O}_P^{\partial\mathbb{B}}(x, z) \leq \partial\delta(x, z) \leq \mathbb{O}^{\partial\mathbb{B}}(x, z) \quad (16b)$$

*Proof.* Here, we show only inequality (16a) since (16b) can be shown similarly. Notice that for  $P = \text{orth}([x, z])$ , the following sequence of inclusions holds:

$$\{PYP^T : Y \in \mathbb{S}^{\text{rank}([x, z])}\} \subseteq \{X \in \mathbb{S}^n : \text{rank}(X) \leq 2r\} \subseteq \mathbb{S}^n. \quad (17)$$

Let  $(\mathcal{H}^*, \delta^*)$  denote the minimizer of the problem corresponding to  $LMI(x, z)$ . By the definition of the  $\mathbb{O}$  function, for every  $X \in \mathbb{S}^n$  it holds that

$$(1 - \delta^*)\|X\|_F^2 \leq \langle X, \mathcal{H}^*(X) \rangle = \|\mathcal{A}^*(X)\|^2 \leq (1 + \delta^*)\|X\|_F^2$$

where the operator  $\mathcal{A}^*$  is such that  $\mathcal{H}^* = \mathcal{A}^{*T} \mathcal{A}^*$  exists because  $\mathcal{H}^* \succeq 0$ . If the inequality holds for all  $X \in \mathbb{S}^n$ , it must hold for  $\text{rank}(X) \leq 2r$ , as noticed by (17). Thus, we conclude that the pair  $(\mathcal{A}^*, \delta^*)$  is feasible for the problem defining  $\delta(x, z)$ . This proves the upper bound. Similarly, if  $(\mathcal{A}_*, \delta_*)$  is the minimizer of the problem defining  $\delta(x, z)$ , then by (17), the pair  $(\mathcal{A}_*^T \mathcal{A}_*, \delta_*)$  is feasible for the problem defining  $\mathbb{O}_P(x, z)$ . This can be verified after rewriting the last constraint of the problem defining  $\mathbb{O}_P$  in the form

$$(1 - \delta)\|PYP^T\|_F^2 \leq \langle PYP^T, \mathcal{A}_*^T \mathcal{A}_*(PYP^T) \rangle = \|\mathcal{A}_*(PYP^T)\|^2 \leq (1 + \delta)\|PYP^T\|_F^2$$

for all  $Y \in \mathbb{S}^{\text{rank}(x, z)}$ . It is important to notice that the same argument works for an arbitrary choice of  $P \in \mathbb{R}^{n \times d}$  with  $d \leq 2r$ .  $\square$

The lemma above completes the proof of Theorem 3. Theorem 2 follows by substituting 1 in the right hand side of (14) and (15). Notice that the linearity of the gradient and the Hessian with respect to the kernel operation matrix is the only property of the objective function that has been extensively used here. It can be exploited for generalization of the developed theory.

Proof of Proposition 4

We write the dual of the problem defining the function  $\mathbb{O}$  as:

$$\begin{aligned} & \underset{y, \lambda, U_1 \geq 0, U_2 \geq 0, V \geq 0}{\text{maximize}} && \text{Tr}[U_1 - U_2] \end{aligned} \quad (18a)$$

$$\text{subject to} \quad \text{Tr}[U_1 + U_2] = 1, \quad (18b)$$

$$\begin{aligned} & \mathcal{L}_{x,z}^T(y) - \mathcal{M}_{x,z}^T(V) + \mathcal{F}^T(\lambda) = \\ & U_1 - U_2 \end{aligned} \quad (18c)$$

This problem is the exact reformulation of

$$\underset{y \in \mathbb{R}^{n \times r}, V \geq 0, \mu \in \mathbb{R}^t}{\text{maximize}} \quad \frac{\sum_{i=1}^d (-\lambda_i (\mathcal{L}_{x,z}^T(y) - \mathcal{M}_{x,z}^T(V) + \mathcal{F}^T(\mu)))_+}{\sum_{i=1}^d (+\lambda_i (\mathcal{L}_{x,z}^T(y) - \mathcal{M}_{x,z}^T(V) + \mathcal{F}^T(\mu)))_+} \quad (19)$$

For details see Lemma 14 by Zhang, Sojoudi, and Lavaei 2019. Both primal and dual problems are bounded and the dual is strictly feasible. Recall vector  $\mathbf{e}$  and matrix  $\mathbf{X}$  from Section 4, such that for all  $u \in \mathbb{R}^{n \times r}$  it holds that

$$\mathbf{e} = \text{vec}(xx^T - zz^T), \quad \mathbf{X}\text{vec}(u) = \text{vec}(xu^T + ux^T)$$

Strictly feasible point of (18) has the components  $y = 0$ ,  $\lambda = 0$ ,  $V = \varepsilon I$ ,  $U_1 = \eta I - \varepsilon W$  and  $U_2 = \eta I + \varepsilon W$  where  $2\eta = n^{-2}$ ,  $2W = r[\text{vec}(I)\mathbf{e}^T + \text{evec}(I)^T] - \mathbf{X}\mathbf{X}^T$ , and  $\varepsilon$  is sufficiently small to ensure that both  $U_1$  and  $U_2$  are PSD. Consequently, Slater's condition and strong duality hold and thus the solution of (19) coincides with  $\mathbb{O}(x, z)$ .

If  $\mathcal{F} = (\mathcal{S}, \mathcal{W})$ , then  $\mathcal{F}^T(u, \mathbf{T}) = \mathcal{W}^T(u) + \mathcal{S}^T(\mathbf{T})$ . At the same time, if  $\mathcal{S}$  is represented by the matrix  $\mathbf{S}$ , then  $\mathcal{S}(\mathbf{T}) = \mathcal{S}^T(\mathbf{T}) = \mathbf{S} \circ \mathbf{T}$ . Let  $\mathbf{S}$  and  $\mathbf{S}'$  be the matrix representations of  $\mathcal{S}$  and  $\mathcal{S}'$ , respectively.  $\mathcal{S}' \subset \mathcal{S}$  means that there exists  $\mathbf{S}^\Delta$  such that  $\mathbf{S} = \mathbf{S}' \cup \mathbf{S}^\Delta$  and  $\mathbf{S}' = \mathbf{S} + \mathbf{S}^\Delta$ . It is straightforward to verify that for any  $\mathbf{R} \in \mathbb{S}^{n^2}$  there exists  $\mathbf{T} \in \mathbb{S}^{n^2}$  such that  $\mathbf{S} \circ \mathbf{T} + \mathbf{S}^\Delta \circ \mathbf{R} = \mathbf{S}' \circ \mathbf{T}$ . The opposite is also true: for any  $\mathbf{T} \in \mathbb{S}^{n^2}$  there exists  $\mathbf{R} \in \mathbb{S}^{n^2}$  such that  $\mathbf{S} \circ \mathbf{T} + \mathbf{S}^\Delta \circ \mathbf{R} = \mathbf{S}' \circ \mathbf{T}$ .

We introduce the short-hand notation  $\mathcal{O}(y, V, u) = \mathcal{L}_{x,z}^T(y) - \mathcal{M}_{x,z}^T(V) + \mathcal{W}^T(u)$ . One can write:

$$\begin{aligned} \mathbb{O}(x, z; \mathcal{W}, \mathcal{S}') &= \underset{y \in \mathbb{R}^{n \times r}, V \geq 0, u \in \mathbb{R}^t, \mathbf{T} \in \mathbb{S}^{n^2}}{\text{minimize}} \quad \frac{\sum_{i=1}^d (-\lambda_i (\mathcal{O}(y, V, u) + \mathbf{S}' \circ \mathbf{T}))_+}{\sum_{i=1}^d (+\lambda_i (\mathcal{O}(y, V, u) + \mathbf{S}' \circ \mathbf{T}))_+} = \\ & \underset{y \in \mathbb{R}^{n \times r}, V \geq 0, u \in \mathbb{R}^t, \mathbf{T} \in \mathbb{S}^{n^2}, \mathbf{R} \in \mathbb{S}^{n^2}}{\text{minimize}} \quad \frac{\sum_{i=1}^d (-\lambda_i (\mathcal{O}(y, V, u) + \mathbf{S} \circ \mathbf{T} + \mathbf{S}^\Delta \circ \mathbf{R}))_+}{\sum_{i=1}^d (+\lambda_i (\mathcal{O}(y, V, u) + \mathbf{S} \circ \mathbf{T} + \mathbf{S}^\Delta \circ \mathbf{R}))_+} \leq \\ & \underset{y \in \mathbb{R}^{n \times r}, V \geq 0, u \in \mathbb{R}^t, \mathbf{T} \in \mathbb{S}^{n^2}}{\text{minimize}} \quad \frac{\sum_{i=1}^d (-\lambda_i (\mathcal{O}(y, V, u) + \mathbf{S} \circ \mathbf{T} + \mathbf{S}^\Delta \circ \mathbf{0}))_+}{\sum_{i=1}^d (+\lambda_i (\mathcal{O}(y, V, u) + \mathbf{S} \circ \mathbf{T} + \mathbf{S}^\Delta \circ \mathbf{0}))_+} = \mathbb{O}(x, z; \mathcal{W}, \mathcal{S}) \end{aligned}$$

This completes the proof.