

LOW RANK APPROXIMATION TO ENTANGLED MULTIPARTITE QUANTUM SYSTEMS

MATTHEW M. LIN* AND MOODY T. CHU†

Abstract. Qualifying the entanglement of a mixed multipartite state by gauging its distance to the nearest separable state of a fixed rank is a challenging but critically important task in quantum technologies. Such a task is computationally demanding partly because of the necessity of optimization over the complex field in order to characterize the underlying quantum properties correctly and partly because of the high nonlinearity due to the multipartite interactions. Representing the quantum states as complex density matrices with respect to some suitably selected bases, this work offers two avenues to tackle this problem numerically. For the rank-1 approximation, an iterative scheme solving a nonlinear singular value problem is investigated. For the general low-rank approximation with probabilistic combination coefficients, a projected gradient dynamics is proposed. Both techniques are shown to converge globally to a local solution. Numerical experiments are carried out to demonstrate the effectiveness and the efficiency of these methods.

Key words. entanglement, separability, multipartite system, low-rank approximation, gradient dynamics, Wirtinger calculus

AMS subject classifications. 65F10, 15A24, 65H10, 15A72, 58D19

1. Introduction. Entanglement is a ubiquitous phenomenon in nature. Given one system, in whatever sense, it is almost inevitable that it will necessarily interact with another or more systems. This involvement of multiple systems, whether tangible or impalpable, can generally be regarded as an entanglement. Depending on the settings, entanglement can be characterized in different forms. Quantum entanglement, where multiple quantum systems interact in such a way as if both their spatial coordinates and their linear momenta are linked, even when the systems are widely separated in space, is particularly intriguing [1]. In modern days, quantum entanglement plays an increasingly more important role in quantum technologies. Quantum informatics and quantum communication, for example, exploit the entanglement for faster and more secure passage of information than classical algorithms. In recent years understanding of entanglement has advanced and diversified into many subfields with applications across a variety of disciplines. The scope is so broad that it is beyond our technical competence, nor is there room in this short note, to provide even the most basic overview of the different subjects related to entanglement. Out of the numerous many, we mention merely three review articles [2, 3, 4] whose references to hundreds of research results on entanglement should be a conspicuous indication of the breadth and the depth of the vast research endeavors in this area. This work concerns only about a fairly focused subject of measuring numerically the distance between a given mixed state and its nearest separable state [5, 6, 7]. In this introduction, therefore, we shall outline only the needed background information pertaining to our methods. For clarity, we divide the discussion into subsections by topics for easy perusal. Readers can skip the parts that they are familiar with.

1.1. Entanglement and separability. In this section we briefly review some basic notion of entanglement and separability. For a more thorough and in-depth treatment of the main ideas, we suggest [8, 9, 10] and the classic book [11].

A quantum mechanical system is typically cast as a complex Hilbert space. The reasons that complex numbers are needed in quantum mechanics are plainly explained in [12, 13] and the references therein. Any unit vector in the space is referred to as a pure state which typically is denoted by the Dirac's ket notation $|\mathbf{x}\rangle$. A mixed quantum state is a probabilistic ensemble of finitely many pure states. It is more convenient to represent a mixed state ϱ as a density matrix

$$\varrho := \sum_i \mu_i |\mathbf{x}_i\rangle \langle \mathbf{x}_i|; \quad \sum_i \mu_i = 1; \quad \mu_i \geq 0, \quad (1.1)$$

*Department of Mathematics, National Cheng Kung University, Tainan 701, Taiwan (mmlin@mail.ncku.edu.tw) This research was supported in part by the National Center for Theoretical Sciences of Taiwan and by the Ministry of Science and Technology of Taiwan under grant 110-2636-M-006 -006.

†Department of Mathematics, North Carolina State University, Raleigh, NC 27695-8205. (chu@math.ncsu.edu) This research was supported in part by the National Science Foundation under grant DMS-1912816.

where the density matrices $|\mathbf{x}_i\rangle\langle\mathbf{x}_i|$ of pure states $|\mathbf{x}_i\rangle$ are simply the orthogonal projector that maps any $|\mathbf{z}\rangle$ onto $|\mathbf{x}_i\rangle\langle\mathbf{x}_i|\mathbf{z}\rangle$ with $\langle\mathbf{x}_i|\mathbf{z}\rangle$ denoting the inner product in the Hilbert space.

A bipartite system $\mathcal{H} = \mathcal{H}_1 \otimes \mathcal{H}_2$ is a composition of two quantum mechanical subsystems \mathcal{H}_1 and \mathcal{H}_2 which interact with each other through a bilinear map¹ denoted by the symbol \otimes . A pure quantum state $|\psi\rangle \in \mathcal{H}$ is called separable if and only if

$$|\psi\rangle = |\psi_1\rangle \otimes |\psi_2\rangle, \quad (1.2)$$

where $|\psi_i\rangle \in \mathcal{H}_i$, $i = 1, 2$, are pure states, respectively; otherwise, $|\psi\rangle$ is entangled. The real issue, however, concerns the mixed state in the composite system. A general (mixed) quantum state $\rho \in \mathcal{H}$ is called separable if it can be decomposed as a probabilistic mixture of tensor products of density matrices of pure states [14, 15]:

$$\rho = \sum_i \theta_i (|\mathbf{x}_i\rangle\langle\mathbf{x}_i|) \otimes (|\mathbf{y}_i\rangle\langle\mathbf{y}_i|), \quad \sum_i \theta_i = 1; \quad \theta_i \geq 0, \quad (1.3)$$

where $|\mathbf{x}_i\rangle \in \mathcal{H}_1$ and $|\mathbf{y}_i\rangle \in \mathcal{H}_2$ are unit vectors. Thus, the collection of all separable states in a bipartite system form a convex set with pure separable states as its extreme points [9].

The same notion of composition can be applied to more than two subsystems [16]. However, the classification of quantum-entangled states is far more complicated than in the bipartite case. On one hand, a natural generalization of (1.3) to a k -partite density matrix ρ is that if

$$\rho = \sum_r \theta_r (|\mathbf{x}_{1,r}\rangle\langle\mathbf{x}_{1,r}|) \otimes \cdots \otimes (|\mathbf{x}_{k,r}\rangle\langle\mathbf{x}_{k,r}|), \quad \sum_i \theta_r = 1; \quad \theta_r \geq 0, \quad (1.4)$$

where, for all r , $|\mathbf{x}_{i,r}\rangle \in \mathcal{H}_i$ is a unit vector, then ρ is called a fully separable state; otherwise, it is said to be fully entangled. On the other hand, there also exists the notion of partially separable states such as the separability with respect to a particular partition of k or the more complicated semi-separability. Once a specific class of separability is chosen, the collection of all separable states under the associated definition still forms a convex set.

Given a general mixed state ρ , if it is not separable, then it is nature to seek its nearest separable state. The task involves calculating the shortest distance between ρ and the convex hull of separable states. This nearest separable approximation problem offers a way to assess the qualification of entanglement. It is of practical importance in quantum applications [2, 11].

1.2. Metric for measurement. We ought to make it clear that the qualification of entanglement depends highly on the assumptions and the applications [17]. For this reason, when measuring the nearness, different metrics might be used for different purposes. We mention three cases.

If the goal is to measure the maximum probability of distinguishability between two quantum states ρ and σ , then the trace metric

$$D_T(\rho, \sigma) := \frac{1}{2} \text{Tr} \sqrt{(\rho - \sigma)^2},$$

based on the so-called Kolmogorov-Smirnov (KS) test for comparing random samples, is perhaps preferred. On the other hand, since repeated measurements are necessary in quantum computation, it might be desired to calculate the minimum number of measurements required to distinguish two different states. For this purpose, the Bures distance

$$D_B(\rho, \sigma) := \sqrt{2 - 2 \text{Tr} \sqrt{\sqrt{\rho} \sigma \sqrt{\rho}}},$$

¹The very same notation \otimes has been used for many different meanings in the literature. The distinction between a tensor product and the Kronecker product is necessary for computation and will be explained in Footnote 2. For a general composite system $\mathcal{H}_1 \otimes \mathcal{H}_2$, we emphasize that \otimes is merely a bilinear map.

an analogue of the Fisher information in classical statistics, can be employed. If we regard the density matrix as an integrated ensemble of the state in which the whole inherent information is contained, then the Frobenius norm

$$D_F(\rho, \sigma) = \frac{1}{2} \|\rho - \sigma\|_F = \frac{1}{2} \sqrt{\text{Tr}(\rho - \sigma)^2}$$

may be used to measure the geometric difference between two ensembles [18].

It is known in linear algebra that, over finite dimensional spaces, all norms are equivalent [19], but in quantum applications different choices of metrics will lead to different approximation results and the associated interpretations. Also, not all distance formulas are easy to use for numerical computation. Taking the positive square root of a positive definite matrix repeatedly in the computation for the metrics D_T or D_B is obviously more expensive than taking the square root of a scalar in the metric D_F . As a starter, we use the Frobenius norm D_F in this work for its ease of implementation. If D_T or D_B is to be used, then specifying the gradient information will be much more involved. It will require separate works to develop new schemes and the pertinent convergence theory. A numerical comparison of various measures is worthy of further investigation, but is beyond the scope of this paper.

1.3. Approximation problem. Suppose \mathcal{H}_1 and \mathcal{H}_2 are two finite dimensional quantum systems with fixed basis states $\{\mathbf{e}_i\}_{i=1}^m$ and $\{\mathbf{f}_j\}_{j=1}^n$, respectively. Then, elements $|\mathbf{x}\rangle \in \mathcal{H}_1$ and $|\mathbf{y}\rangle \in \mathcal{H}_2$ can be interpreted as two column vectors $\mathbf{x} \in \mathbb{C}^m$ and $\mathbf{y} \in \mathbb{C}^n$ of their coordinates, respectively. The density matrices $|\mathbf{x}\rangle\langle\mathbf{x}|$ and $|\mathbf{y}\rangle\langle\mathbf{y}|$ are indeed rank-1 matrices with unit trace in $\mathbb{C}^{m \times m}$ and $\mathbb{C}^{n \times n}$, respectively. Furthermore, with respect to the basis $\mathbf{e}_i \otimes \mathbf{f}_j$ in the lexicographical order, the tensor product can be interpreted as the Kronecker product. Therefore, the approximation problem

$$\min_{\substack{|\mathbf{x}_i\rangle \in \mathcal{H}_1, \langle\mathbf{x}_i|\mathbf{x}_i\rangle=1 \\ |\mathbf{y}_i\rangle \in \mathcal{H}_2, \langle\mathbf{y}_i|\mathbf{y}_i\rangle=1 \\ \sum_i \theta_i=1, \theta_i \geq 0}} \|\rho - \sum_i \theta_i (|\mathbf{x}_i\rangle\langle\mathbf{x}_i|) \otimes (|\mathbf{y}_i\rangle\langle\mathbf{y}_i|)\|_F^2, \quad (1.5)$$

can be translated via the linear algebra interpretation into the following equivalent problem

$$\min_{\substack{\mathbf{x}_r \in \mathbb{C}^m, \|\mathbf{x}_r\|=1, \\ \mathbf{y}_r \in \mathbb{C}^n, \|\mathbf{y}_r\|=1, \\ \lambda_r \geq 0, \sum_r \lambda_r=1}} \|\rho - \sum_{r=1}^R \lambda_r (\mathbf{x}_r \mathbf{x}_r^*) \otimes (\mathbf{y}_r \mathbf{y}_r^*)\|_F^2, \quad (1.6)$$

where $\rho \in \mathbb{C}^{mn \times mn}$ is positive definite (hence hermitian) with unit trace, $*$ denotes the conjugate transpose, and \otimes is interpreted as the Kronecker product. Over the framework of general Hilbert spaces, the term needed for the summation in (1.5) is difficult to determine. Over the finite dimensional spaces we know by the Carathéodory theorem [20, Theorem 2.2.4] that no more than $(mn)^2 + 1$ terms will provide the best approximation of ρ over the convex hull of separable states. The problem therefore involves at most $(2(m+n)+1)((mn)^2+1)$ real variables. Suppose that R is a predetermined positive integer, then we have a low-rank approximation problem. In this case, since we are not taking all the extreme points of the convex hull of the pure states into the summation, the solution to (1.6) is not unique.

This paper concerns the general k -partite low-rank approximation problem of the form

$$\min_{\substack{\mathbf{x}_{i,r} \in \mathbb{C}^{m_i}, \|\mathbf{x}_{i,r}\|_2=1 \\ \lambda_r \geq 0, \sum_{r=1}^R \lambda_r=1}} \|\rho - \sum_{r=1}^R \lambda_r (\mathbf{x}_{1,r} \mathbf{x}_{1,r}^*) \otimes \cdots \otimes (\mathbf{x}_{k,r} \mathbf{x}_{k,r}^*)\|_F^2, \quad (1.7)$$

for a given density matrix $\rho \in \mathbb{C}^{\prod_{i=1}^k m_i \times \prod_{i=1}^k m_i}$ and $k \geq 2$. It might appear that we are dealing with the full separability for a k -partite system. Nevertheless, our techniques applied to general dimensions m_i . It is possible that a single space \mathbb{C}^{m_i} contains the composition of several subsystems. That is, our methods can be applied to explore the partial separability approximation as well [21]. This can be best illustrated by the split of an n -qubit system in the next section.

1.4. Qubit system. The setting we present in this work is over a general multipartite quantum mechanical system with $\mathbf{x}_{i,r} \in \mathbb{C}^{m_i}$, where m_i is an arbitrary positive integer. For applications in quantum information science, a commonly used basic unit for quantum computation is the 2-dimensional Hilbert space \mathbb{C}^2 . In this context, we still can formulate the low-rank approximation.

Denoting the canonical basis vectors over \mathbb{C}^2 denoted by $|0\rangle = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ and $|1\rangle = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ or simply $|\uparrow\rangle$ and $|\downarrow\rangle$, a qubit is the quantum mechanical analogue of a classical bit in the digital computer. Correspondingly, in the bipartite system $\mathbb{C}^2 \otimes \mathbb{C}^2$ the product $|\uparrow\rangle \otimes |\downarrow\rangle$ is often abbreviated as $|\uparrow\downarrow\rangle$, referred to as a 2-qubit. A d -qubit system is represented by $(\mathbb{C}^2)^{\otimes d} = \mathbb{C}^2 \otimes \dots \otimes \mathbb{C}^2$. Therefore, a state in the system can be thought of as a complex vector of dimension 2^d . One could regard $(\mathbb{C}^2)^{\otimes d}$ as an d -partite entangled system of \mathbb{C}^2 . If we regard the zeros and ones as constituting the binary expansion of an integer, say, ℓ , then we can replace the representations of a basic d -qubit state by a short form $|\ell\rangle$, $0 \leq \ell < 2^d - 1$. On the other hand, if we split $d = p + q$, then we could also consider $(\mathbb{C}^2)^{\otimes d} = (\mathbb{C}^2)^{\otimes p} \otimes (\mathbb{C}^2)^{\otimes q}$ as a bipartite entanglement of $(\mathbb{C}^2)^{\otimes p}$ and $(\mathbb{C}^2)^{\otimes q}$. In the latter case, the problem (1.5) becomes a partial separability approximation with $m = 2^p$ and $n = 2^q$. Given a 32×32 density matrix ρ , for example, we can group the 5 qubits in 7 ways: $5 = 5 + 0 = 4 + 1 = 3 + 2 = 3 + 1 + 1 = 2 + 2 + 1 = 2 + 1 + 1 + 1 = 1 + 1 + 1 + 1 + 1$, each constitutes a distinct low-rank approximation problems. The techniques to be described in this paper can be applied to handle each case with appropriate realization of k and m_i in (1.7).

1.5. Canonical polyadic decomposition. Before we move on to describe our numerical method, we ought to point out that, for the case $R = 1$, the problem (1.7) can be recast as a specially structured low-rank tensor approximation referred to as the canonical polyadic decomposition with symmetry in the literature [22]. For example, if $\rho \in \mathbb{R}^{mn \times mn}$ is properly folded into an order-4 tensor $\mathfrak{A} \in \mathbb{R}^{m \times m \times n \times n}$, we may recast the real version of (1.6) as an order-4 rank-1 tensor approximation with symmetry in the first two and the last two modes:

$$\min_{\substack{\lambda \in \mathbb{R}_+, \mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n \\ \|\mathbf{x}\|=1, \|\mathbf{y}\|=1}} \|\mathfrak{A} - \lambda \mathbf{x} \circ \mathbf{x} \circ \mathbf{y} \circ \mathbf{y}\|_F^2, \quad (1.8)$$

where \circ denotes the outer product. Many techniques, e.g., those in the `Tensorlab` toolbox [23], are readily available to handle this specially structured rank-1 problem. For the case $R > 1$, however, it becomes challenging to satisfy the probabilistic constraint by conventional techniques. So far as we know, the `Tensorlab` toolbox has not developed this functionality yet. Recall that the probabilistic ensemble is essential in quantum applications. One of our contributions in this work is a mechanism to maintain this constraint.

This paper is organized as follows. In Section 2, we generalize our recent results for real-valued bipartite rank-1 approximation [24] to complex-valued multipartite rank-1 approximation. This generalization prepares the way of using the Wirtinger calculus to derive the gradient of a real-valued objective function with complex variables. In order to address the probabilistic constraint, we propose in Section 3 a projected gradient flow to tackle the multipartite low-rank approximation (1.7) directly. The most important features of this dynamical system are that the nonnegativity and sum-to-one constraints of the combinations coefficients are preserved and that the rank can be automatically adjusted downward during the integration. We believe that the simplicity of this approach might be employed as a useful tool entanglement qualification. Numerical experiments are carried out in Section 4 to demonstrate the working of our algorithms.

2. Multipartite Rank-1 Approximation. In an earlier study [24], we have considered the problem of approximating a real-valued, symmetric and positive matrix $A \in \mathbb{R}^{mn \times mn}$ by a real-valued rank-1 bipartite system, i.e,

$$\min_{\mathbf{x} \in S^{m-1}, \mathbf{y} \in S^{n-1}, \lambda \in \mathbb{R}_+} \|A - \lambda(\mathbf{x}\mathbf{x}^\top) \otimes (\mathbf{y}\mathbf{y}^\top)\|_F^2. \quad (2.1)$$

The idea is to reformulate (2.1) as either a nonlinear eigenvalue problem or a nonlinear singular value problem. Correspondingly, a nonlinear power-like and a nonlinear SVD-like iterative schemes have been proposed and

analyzed. Numerical experiments suggest that these methods are not only easy to implement, but also are highly efficient when comparing with the more sophisticated routines used in the package Tensorlab.

In this section, we consider the most general k -partite approximation problem

$$\min_{\substack{\lambda \in \mathbb{R}_+, i \in \llbracket k \rrbracket \\ \mathbf{x}_i \in \mathbb{C}^{m_i}, \|\mathbf{x}_i\|_2=1}} \|A - \lambda(\mathbf{x}_1 \mathbf{x}_1^* \otimes \cdots \otimes \mathbf{x}_k \mathbf{x}_k^*)\|_F^2, \quad (2.2)$$

for a density matrix $A \in \mathbb{C}^{\prod_{i=1}^k m_i \times \prod_{i=1}^k m_i}$. The previous experiences we have learned in [24] help but the generalization is not as obvious because of the involvement of complex variables and the extended number of factors when $k \geq 3$.

2.1. Basics. To facilitate our subsequent discussion, we first introduce some basic notations and review some useful facts. Let $\llbracket k \rrbracket$ denote the set of integers $\{1, \dots, k\}$. Given column vectors $\mathbf{x}_i, i \in \llbracket k \rrbracket$, note that the classical Kronecker product \otimes is equivalent to the tensor product² \circ in a reversed order [25], i.e.,

$$\mathbf{x}_1 \otimes \cdots \otimes \mathbf{x}_k = \text{vec}(\mathbf{x}_k \circ \cdots \circ \mathbf{x}_1), \quad (2.3)$$

where $\text{vec}(\mathfrak{T})$ for an order- k tensor $\mathfrak{T} \in \mathbb{C}^{m_1 \times m_2 \times \cdots \times m_k}$ is a linear array whose entry at the location

$$(s_k - 1)m_{k-1} \dots m_1 + (s_{k-1} - 1)m_{k-2} \dots m_1 + \dots + (s_2 - 1)m_1 + s_1$$

is precisely the element τ_{s_1, \dots, s_k} of \mathfrak{T} . It will be convenient to adopt the abbreviations

$$\begin{cases} \bigotimes_{i=1}^k \mathbf{x}_i & := \mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \cdots \otimes \mathbf{x}_k, \\ \bigcirc_{i=k}^1 \mathbf{x}_i & := \mathbf{x}_k \circ \cdots \circ \mathbf{x}_1, \end{cases}$$

and define the order- k tensor

$$\mathfrak{D}(\mathbf{x}_1, \dots, \mathbf{x}_k) := \text{reshape}(A \bigotimes_{i=1}^k \mathbf{x}_i, [m_k, \dots, m_1]) \in \mathbb{C}^{m_k \times \cdots \times m_1},$$

where the operator `reshape` is identical to that in `Matlab` which returns a multi-dimensional array with the specified dimensions.

To handle the multi-indices more effectively, the following notation system proves handy [26]. Suppose that the set $\llbracket k \rrbracket$ is partitioned as the union of two disjoint nonempty subsets $\alpha := \{\alpha_1, \dots, \alpha_\ell\}$ and $\beta := \{\beta_1, \dots, \beta_{k-\ell}\}$. Let $\mathcal{I} = (i_1, \dots, i_\ell)$ and $\mathcal{J} = (j_1, \dots, j_{k-\ell})$ denote indices at locations α and β , respectively, where each index in the arrays \mathcal{I} and \mathcal{J} should be within the corresponding range of integers, e.g., $i_1 \in \llbracket m_{\alpha_1} \rrbracket$ and so on. An element $\tau_{s_1 \dots s_k}$ in the order- k tensor \mathfrak{T} can be identified as $\tau_{[\mathcal{I}|\mathcal{J}]}$ with $s_{\alpha_\mu} = i_\mu$ and $s_{\beta_\nu} = j_\nu$, $\mu \in \llbracket \ell \rrbracket$, $\nu \in \llbracket k - \ell \rrbracket$. The point to make is that via the location pointer (α, β) we can enumerate elements $\tau_{s_1 \dots s_k}$ in any order we want. When the reference to a specific partitioning (α, β) is clear, we abbreviate the element as $\tau_{[\mathcal{I}|\mathcal{J}]}$. The partition (α, β) may be regarded as generalizing the familiar notion of rows and columns for matrices.

Given a partition $\llbracket k \rrbracket = \alpha \cup \beta$, we may regard an order- k tensor \mathfrak{T} as the matrix representation of a linear transformation from the tensor space $\mathbb{C}^{m_{\beta_1} \times \cdots \times m_{\beta_{k-\ell}}}$ to $\mathbb{C}^{m_{\alpha_1} \times \cdots \times m_{\alpha_\ell}}$. Thus, we use the symbol \circledast_α to replace the conventional "matrix-to-vector" multiplication, that is, if $\mathfrak{S} = [\sigma_{j_1 \dots j_{k-\ell}}] \in \mathbb{C}^{m_{\beta_1} \times \cdots \times m_{\beta_{k-\ell}}}$, then the \mathcal{I} -th element of the product $\mathfrak{Q} = \mathfrak{T} \circledast_\alpha \mathfrak{S} \in \mathbb{C}^{m_{\alpha_1} \times \cdots \times m_{\alpha_\ell}}$ is defined by

$$(\mathfrak{Q})_{\mathcal{I}} := \sum_{j_1=1}^{m_{\beta_1}} \cdots \sum_{j_{k-\ell}=1}^{m_{\beta_{k-\ell}}} \tau_{[\mathcal{I}|j_1 \dots j_{k-\ell}]}^{(\alpha, \beta)} \sigma_{j_1, \dots, j_{k-\ell}}.$$

² The tensor product of tensors leads to a multi-indexed array. While the way to enumerate its elements is often immaterial in theory, it is essential to enumerate them consistently for numerical calculation. One general rule adopted is that the indices of the leftmost tensor are counted first, e.g., the indices in the tensor product $\mathbf{a} \circ \mathbf{b}$ of two vectors are enumerated in the same way as the matrix $\mathbf{a}\mathbf{b}^\top$. The relationship (2.3) therefore follows.

For $\mathbf{a}, \mathbf{b} \in \mathbb{C}^n$, let

$$\langle \mathbf{a}, \mathbf{b} \rangle_{\mathbb{R}} := \sum_{i=1}^n a_i b_i \quad (2.4)$$

denote a formal inner product. Similarly, the notation $\langle \cdot, \cdot \rangle_{\mathbb{R}}$ can be generalized to matrices or tensors. The relationship

$$\langle \mathfrak{T}, \bigcirc_{i=1}^k \mathbf{x}_i \rangle_{\mathbb{R}} = \langle \mathfrak{T} \otimes_{\alpha} \left(\bigcirc_{s=1}^{k-\ell} \mathbf{x}_{\beta_s} \right), \bigcirc_{t=1}^{\ell} \mathbf{x}_{\alpha_t} \rangle_{\mathbb{R}}, \quad (2.5)$$

which is nothing but the associative law of multiplication, holds for any tensors $\bigcirc_{t=1}^{\ell} \mathbf{x}_{\alpha_t} \in \mathbb{C}^{m_{\alpha_1} \times \dots \times m_{\alpha_{\ell}}}$ and $\bigcirc_{s=1}^{k-\ell} \mathbf{x}_{\beta_s} \in \mathbb{C}^{m_{\beta_1} \times \dots \times m_{\beta_{k-\ell}}}$. We shall employ (2.5) to help describe lengthy algebraic manipulations.

Suppose that $f : \mathbb{C} \rightarrow \mathbb{R}$ is a real-valued function over a complex variable $z = x + iy$. If we regard $f(z) = u(x, y)$, then the Wirtinger derivatives are defined by

$$\begin{cases} \frac{\partial f}{\partial z} & := \frac{1}{2} \left(\frac{\partial u}{\partial x} - i \frac{\partial u}{\partial y} \right), \\ \frac{\partial f}{\partial \bar{z}} & := \frac{1}{2} \left(\frac{\partial u}{\partial x} + i \frac{\partial u}{\partial y} \right). \end{cases}$$

That is, while maintaining the usual complex arithmetic throughout the operations, we take the formal partial derivatives of $f(z)$ by treating z and \bar{z} as independent variables with respect to each other [27]. For a general real-valued function $f : \mathbb{C}^n \rightarrow \mathbb{R}$, the definition of the Wirtinger derivative can be generalized to:

$$\begin{cases} \frac{\partial f}{\partial z} & := \frac{1}{2} \left(\frac{\partial f}{\partial \mathbf{u}} - i \frac{\partial f}{\partial \mathbf{v}} \right), \\ \frac{\partial f}{\partial \bar{z}} & := \frac{1}{2} \left(\frac{\partial f}{\partial \mathbf{u}} + i \frac{\partial f}{\partial \mathbf{v}} \right), \end{cases} \quad (2.6)$$

where we regard $f(\mathbf{z}) = f(\mathbf{u}, \mathbf{v})$ in the real variables $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ and $\mathbf{z} = \mathbf{u} + i\mathbf{v} \in \mathbb{C}^n$. In this way, the ‘‘true’’ gradient of function $f : \mathbb{C}^n \rightarrow \mathbb{R}$ can be calculated from the Wirtinger derivatives via the relationship:

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial \mathbf{u}} \\ \frac{\partial f}{\partial \mathbf{v}} \end{bmatrix} = \begin{bmatrix} \frac{\partial f}{\partial \mathbf{z}} + \frac{\partial f}{\partial \bar{\mathbf{z}}} \\ i \left(\frac{\partial f}{\partial \mathbf{z}} - \frac{\partial f}{\partial \bar{\mathbf{z}}} \right) \end{bmatrix}. \quad (2.7)$$

2.2. Nonlinear Singular Value Formulation. For the optimization problem (2.2), and especially for the case $k \geq 3$, we propose the idea of alternately applying the singular value decomposition to update two complex vectors at a time. We divide the discussion into two parts. First, we motivate the iterative scheme by exploring the first order optimal condition for the objective function. Then, we derive the convergence theory.

LEMMA 2.1. *Let $\llbracket k \rrbracket = \alpha \cup \beta$ with $\alpha := \{\alpha_1, \alpha_2\}$ and $\beta := \{\beta_1, \dots, \beta_{k-2}\}$ be an arbitrary partition. If $(\mathbf{x}_1, \dots, \mathbf{x}_k)$ is a local minimizer to (2.2), then it is necessary that*

$$\begin{cases} (\mathfrak{D}(\mathbf{x}_1, \dots, \mathbf{x}_k) \otimes_{\alpha} \left(\bigcirc_{i=1}^{k-2} \bar{\mathbf{x}}_{\beta_i} \right) \bar{\mathbf{x}}_{\alpha_2}) & = \lambda(\mathbf{x}_1, \dots, \mathbf{x}_k) \mathbf{x}_{\alpha_1}, \\ (\mathfrak{D}(\mathbf{x}_1, \dots, \mathbf{x}_k) \otimes_{\alpha} \left(\bigcirc_{i=1}^{k-2} \bar{\mathbf{x}}_{\beta_i} \right) \mathbf{x}_{\alpha_1})^* & = \lambda(\mathbf{x}_1, \dots, \mathbf{x}_k) \bar{\mathbf{x}}_{\alpha_2}. \end{cases} \quad (2.8)$$

It is worth noting that the multiplication between the order- k tensor $\mathfrak{D}(\mathbf{x}_1, \dots, \mathbf{x}_k)$ and the order-($k-2$) tensor $\bigcirc_{i=1}^{k-2} \bar{\mathbf{x}}_{\beta_i}$ results a matrix. Also, since $\llbracket k \rrbracket = \alpha \cup \beta$ is an arbitrary partition, the specifics of α_1 and α_2 are immaterial. They refer to every possible indices. The necessary condition (2.8) therefore is much more involved than it appears.

Proof. Because A is positive definite, the minimization of (2.2) is equivalent to the maximization of

$$\max_{\substack{\mathbf{x}_i \in \mathbb{C}^{m_i}, \|\mathbf{x}_i\|_2=1 \\ i \in \llbracket k \rrbracket}} \lambda(\mathbf{x}_1, \dots, \mathbf{x}_k) := \langle A, (\mathbf{x}_1 \otimes \dots \otimes \mathbf{x}_k)(\mathbf{x}_1 \otimes \dots \otimes \mathbf{x}_k)^* \rangle, \quad (2.9)$$

where $\langle \cdot, \cdot \rangle$ denotes the Frobenius inner product over the complex space. We can also write λ as

$$\lambda(\mathbf{x}_1, \dots, \mathbf{x}_k) = \langle A, \overline{(\mathbf{x}_1 \otimes \dots \otimes \mathbf{x}_k)}(\mathbf{x}_1 \otimes \dots \otimes \mathbf{x}_k)^\top \rangle_{\mathbb{R}}.$$

Consider the variable $\mathbf{x}_{\alpha_1} = \mathbf{u}_{\alpha_1} + i\mathbf{v}_{\alpha_1}$ first. Taking the Wirtinger derivatives with respect to $\overline{\mathbf{x}}_{\alpha_1}$ yields

$$\frac{\partial \lambda}{\partial \overline{\mathbf{x}}_{\alpha_1}} = (\mathfrak{D}(\mathbf{x}_1, \dots, \mathbf{x}_k) \otimes_{\alpha} \left(\bigcirc_{i=1}^{k-2} \overline{\mathbf{x}}_{\beta_i} \right)) \overline{\mathbf{x}}_{\alpha_2}$$

Since

$$\lambda(\mathbf{x}_1, \dots, \mathbf{x}_k) = \overline{\lambda(\mathbf{x}_1, \dots, \mathbf{x}_k)} = \langle \overline{A}, (\mathbf{x}_1 \otimes \dots \otimes \mathbf{x}_k)(\overline{\mathbf{x}}_1 \otimes \dots \otimes \overline{\mathbf{x}}_k)^\top \rangle_{\mathbb{R}},$$

we also have

$$\frac{\partial \lambda}{\partial \mathbf{x}_{\alpha_1}} = \overline{\left(\frac{\partial \lambda}{\partial \overline{\mathbf{x}}_{\alpha_1}} \right)} = \overline{(\mathfrak{D}(\mathbf{x}_1, \dots, \mathbf{x}_k) \otimes_{\alpha} \left(\bigcirc_{i=1}^{k-2} \overline{\mathbf{x}}_{\beta_i} \right)) \overline{\mathbf{x}}_{\alpha_2}},$$

It follows from (2.7) that the partial gradient of λ with respect to the real variables \mathbf{u}_{α_1} and \mathbf{v}_{α_1} is given by

$$\nabla_{(\mathbf{u}_{\alpha_1}, \mathbf{v}_{\alpha_1})} \lambda = \begin{bmatrix} \frac{\partial \lambda}{\partial \mathbf{u}_{\alpha_1}} \\ \frac{\partial \lambda}{\partial \mathbf{v}_{\alpha_1}} \end{bmatrix} = 2 \begin{bmatrix} \mathcal{R} \\ \mathcal{I} \end{bmatrix}, \quad (2.10)$$

where \mathcal{R} and \mathcal{I} are, respectively, the real and imaginary parts of

$$\mathfrak{D}(\mathbf{x}_1, \dots, \mathbf{x}_k) \otimes_{\alpha} \left(\bigcirc_{i=1}^{k-2} \overline{\mathbf{x}}_{\beta_i} \right) \overline{\mathbf{x}}_{\alpha_2} = \mathcal{R} + i\mathcal{I}.$$

Let $S^{2m_{\alpha_1}-1}$ denote the unit sphere

$$S^{2m_{\alpha_1}-1} := \left\{ \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} \in \mathbb{R}^{2m_{\alpha_1}} \mid \|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2 = 1 \right\}.$$

The projection of $\nabla_{(\mathbf{u}_{\alpha_1}, \mathbf{v}_{\alpha_1})} \lambda$ onto the unit sphere S^{2m_j-1} is given by

$$2 \begin{bmatrix} \mathcal{R} \\ \mathcal{I} \end{bmatrix} - 2(\mathbf{u}_{\alpha_1}^* \mathcal{R} + \mathbf{v}_{\alpha_1}^* \mathcal{I}) \begin{bmatrix} \mathbf{u}_{\alpha_1} \\ \mathbf{v}_{\alpha_1} \end{bmatrix}. \quad (2.11)$$

Observe that

$$\lambda = \bar{\lambda} = (\mathbf{u}_{\alpha_1}^* - i\mathbf{v}_{\alpha_1}^*)(\mathcal{R} + i\mathcal{I}) = (\mathbf{u}_{\alpha_1}^* \mathcal{R} + \mathbf{v}_{\alpha_1}^* \mathcal{I}) + i(\mathbf{u}_{\alpha_1}^* \mathcal{I} - \mathbf{v}_{\alpha_1}^* \mathcal{R}).$$

Therefore it must be that

$$\begin{cases} \mathbf{u}_{\alpha_1}^\top \mathcal{R} + \mathbf{v}_{\alpha_1}^\top \mathcal{I} = \lambda, \\ \mathbf{u}_{\alpha_1}^\top \mathcal{I} - \mathbf{v}_{\alpha_1}^\top \mathcal{R} = 0. \end{cases} \quad (2.12)$$

Algorithm 1 (Best rank-1 approximation via SVD updating with randomization.)

Require: An density matrix A and k starting unit vectors $\mathbf{x}_i^{[0]} \in \mathbb{C}^{m_i}$, $i \in \llbracket k \rrbracket$.

Ensure: A local best rank-1 approximation to A in the sense of (2.2).

```

1:  $p \leftarrow 0$ 
2:  $\lambda^{[0]} \leftarrow \langle A, (\bigotimes_{i=1}^k \mathbf{x}_i^{[0]})(\bigotimes_{i=1}^k \mathbf{x}_i^{[0]})^* \rangle$ 
3: repeat
4:    $p \leftarrow p + 1$ 
5:    $\alpha \leftarrow$  two random integers from  $\llbracket k \rrbracket$ 
6:    $\beta \leftarrow \llbracket k \rrbracket - \alpha$  { Complement of  $\alpha$  }
7:    $D \leftarrow \mathfrak{D}(\mathbf{x}_1^{[p-1]}, \dots, \mathbf{x}_k^{[p-1]}) \circledast_{\alpha} \left( \bigcirc_{i=1}^{k-2} \bar{\mathbf{x}}_{\beta_i}^{[p-1]} \right)$ 
8:    $[\mathbf{u}, s, \mathbf{v}] = \text{svds}(D, 1)$  { Dominant singular value triplets via, e.g., Matlab routine svds }
9:    $\theta \leftarrow$  argument of first entry of  $\mathbf{u}$ 
10:   $\mathbf{x}_{\alpha_1}^{[p]} = e^{-i\theta} \mathbf{u}$ 
11:   $\mathbf{x}_{\alpha_2}^{[p]} = e^{i\theta} \bar{\mathbf{v}}$ 
12:   $\lambda^{[p]} \leftarrow s$ 
13: until  $\lambda^{[p]}$  meets convergence criteria

```

The first order optimality condition requires that the projected gradient in any direction be zero. By substituting (2.12) into (2.11), we find that

$$\left(\mathfrak{D}(\mathbf{x}_1 \otimes \dots \otimes \mathbf{x}_k) \circledast_{\alpha} \left(\bigcirc_{i=1}^{k-2} \bar{\mathbf{x}}_{\beta_i} \right) \right) \bar{\mathbf{x}}_{\alpha_2} = \lambda(\mathbf{x}_1, \dots, \mathbf{x}_k) \mathbf{x}_{\alpha_1}.$$

which is the first equation in (2.8). The second equation can be proved by applying a similar argument to the variable $\bar{\mathbf{x}}_{\alpha_2}$. \square

Since the goal is to maximize $\lambda(\mathbf{x}_1, \dots, \mathbf{x}_k)$, we can interpret the relationship (2.8) in Lemma 2.1 in terms of the singular value decomposition as follows.

COROLLARY 2.2. *With respect to an arbitrary but fixed partition $\llbracket k \rrbracket = \alpha \cup \beta$ with $\alpha := \{\alpha_1, \alpha_2\}$ and $\beta := \{\beta_1, \dots, \beta_{k-2}\}$, the triplets $(\mathbf{x}_{\alpha_1}, \lambda, \bar{\mathbf{x}}_{\alpha_2})$ such that (2.8) is satisfied and such that λ is as large as possible must be the dominant singular triplets of the matrix $\mathfrak{D}(\mathbf{x}_1, \dots, \mathbf{x}_k) \circledast_{\alpha} \left(\bigcirc_{i=1}^{k-2} \bar{\mathbf{x}}_{\beta_i} \right)$. In particular, \mathbf{x}_{α_1} is the dominant left singular vector and $\bar{\mathbf{x}}_{\alpha_2}$ is the dominant right singular vector.*

Corollary 2.2 thus motivates an SVD-like iteration where we update two pure states at a time by varying the indices in α . The selections of α could be systematic such as cycling through the list of pairs (1, 2), (2, 3), \dots , $(k-1, k)$ and $(k, 1)$, or could be randomly generated at every iteration. Our proof of convergence does not depend on how α is generated. The updating scheme with random selection of α is sketched in Algorithm 1.

A general purpose routine, say, `svds`, is employed as a black box to calculate the dominant singular triplets. To ensure continuity, we shall align all singular vectors by requiring that the first entries of left singular vectors be real and nonnegative. This can easily be accomplished by a phase change. For example, if $(\mathbf{u}, s, \mathbf{v})$ represents the dominant singular triplets of a matrix X , i.e.,

$$X\mathbf{v} = s\mathbf{u},$$

then so does the triplets $(e^{-i\theta}\mathbf{u}, s, e^{-i\theta}\mathbf{v})$ for any angle θ . Taking θ to be the phase of the first entry of \mathbf{u} will make the first entry of $e^{-i\theta}\mathbf{u}$ nonnegative. This mechanism is included in Algorithm 1.

For the sake of conveniently registering the iterates for analysis, we have implied in the description of Algorithm 1 that whenever two vectors $(\mathbf{x}_{\alpha_1}^{[p]}, \mathbf{x}_{\alpha_2}^{[p]})$ are updated to $(\mathbf{x}_{\alpha_1}^{[p+1]}, \mathbf{x}_{\alpha_2}^{[p+1]})$, the remaining list in $(\mathbf{x}_1^{[p+1]}, \dots, \mathbf{x}_k^{[p+1]})$ are just exact copies of $(\mathbf{x}_{\beta_1}^{[p]}, \dots, \mathbf{x}_{\beta_{k-2}}^{[p]})$, i.e., $\mathbf{x}_{\beta_i}^{[p+1]} = \mathbf{x}_{\beta_i}^{[p]}$ for $i \in \llbracket k-2 \rrbracket$.

In our current application, $\mathfrak{D}(\mathbf{x}_1^{[p]}, \dots, \mathbf{x}_k^{[p]})$ varies in p . This is different from the SVD-based methods developed earlier for stationary tensor approximations [26, 28]. We have to establish a new convergence theory. Toward that goal, we first prove the monotone increasing property of the objective values of λ .

THEOREM 2.3. *Given a density matrix $A \in \mathbb{C}^{mn \times mn}$, let $\{\lambda^{[p]}\}$ be the sequence generated by Algorithm 1 where $\alpha \in \llbracket k \rrbracket$ is randomly selected. Then the inequalities*

$$\lambda(\mathbf{x}_1^{[p]}, \dots, \mathbf{x}_k^{[p]}) \leq \lambda^{[p+1]} \leq \lambda(\mathbf{x}_1^{[p+1]}, \dots, \mathbf{x}_k^{[p+1]}) \leq \lambda^{[p+2]} \quad (2.13)$$

hold. Therefore, both sequences $\{\lambda^{[p]}\}$ and $\{\lambda(\mathbf{x}_1^{[p]}, \dots, \mathbf{x}_k^{[p]})\}$ converge monotonically.

Proof. Define the abbreviation $\mathbf{a}^{[p]} = \mathbf{x}_1^{[p]} \otimes \dots \otimes \mathbf{x}_k^{[p]}$. Then we can write

$$\begin{aligned} \lambda(\mathbf{x}_1^{[p]}, \dots, \mathbf{x}_k^{[p]}) &= \langle \mathbf{A}\mathbf{a}^{[p]}, \mathbf{a}^{[p]} \rangle = \langle \mathfrak{D}(\mathbf{x}_1^{[p]}, \dots, \mathbf{x}_k^{[p]}) \otimes_{\alpha} \left(\bigcirc_{i=1}^{k-2} \overline{\mathbf{x}}_{\beta_i}^{[p]} \right), \bigcirc_{i=1}^2 \mathbf{x}_{\alpha_i}^{[p]} \rangle, \\ \lambda^{[p+1]} &= \langle \mathbf{A}\mathbf{a}^{[p+1]}, \mathbf{a}^{[p+1]} \rangle = \langle \mathfrak{D}(\mathbf{x}_1^{[p]}, \dots, \mathbf{x}_k^{[p]}) \otimes_{\alpha} \left(\bigcirc_{i=1}^{k-2} \overline{\mathbf{x}}_{\beta_i}^{[p]} \right), \bigcirc_{i=1}^2 \mathbf{x}_{\alpha_i}^{[p+1]} \rangle. \end{aligned}$$

The first inequality in (2.13) follows from the definition that $\lambda^{[p+1]}$ is the dominant singular value of the matrix $\mathfrak{D}(\mathbf{x}_1^{[p]}, \dots, \mathbf{x}_k^{[p]}) \otimes_{\alpha} \left(\bigcirc_{i=1}^{k-2} \overline{\mathbf{x}}_{\beta_i}^{[p]} \right)$. Similarly, the third inequality holds. To prove the second inequality, observe that

$$\begin{aligned} \lambda(\mathbf{a}^{[p+1]}) - \lambda^{[p+1]} &= \langle \mathbf{A}\mathbf{a}^{[p+1]}, \mathbf{a}^{[p+1]} \rangle - \langle \mathbf{A}\mathbf{a}^{[p]}, \mathbf{a}^{[p+1]} \rangle = \langle \mathbf{a}^{[p+1]} - \mathbf{a}^{[p]}, \mathbf{A}\mathbf{a}^{[p+1]} \rangle \\ &= \langle \mathbf{a}^{[p+1]} - \mathbf{a}^{[p]}, \mathbf{A}(\mathbf{a}^{[p+1]} - \mathbf{a}^{[p]}) \rangle + \langle \mathbf{a}^{[p+1]} - \mathbf{a}^{[p]}, \mathbf{A}\mathbf{a}^{[p]} \rangle \geq 0, \end{aligned}$$

which completes the proof. \square

We next prove the convergence of iterates themselves under the following generic condition.

DEFINITION 2.4. *We say that the matrix A satisfies Condition P if the corresponding polynomial system (2.8) has finitely many geometrically isolated real-valued solutions.*

Though pathological examples can be constructed, it is well known in algebraic geometry that almost every square system of polynomial equations over the complex field has finitely many solutions [29]. Furthermore, if $F(\mathbf{z}; \mathbf{q})$ is a system of polynomials in both the variables \mathbf{z} and the parameters \mathbf{q} , and is square in \mathbf{z} , then for almost all parameters \mathbf{q} the number of geometrically isolated solutions to this polynomial system is finite [30, Theorem 7.1.1]. The phrase ‘‘almost all’’ means that those values of parameters that fail to produce finitely many and geometrically isolated solutions constitute a nowhere dense and measure zero subset in the ambient space. These exceptions are referred to as ‘‘non-generic’’. For this reason, the condition P is generic.

THEOREM 2.5. *Suppose that the given density matrix $A \in \mathbb{C}^{mn \times mn}$ satisfies the Condition P. Suppose also that the matrices $\mathfrak{D}(\mathbf{x}_1^{[p]}, \dots, \mathbf{x}_k^{[p]})$ always have simple dominant singular values. Then the corresponding iterates $\{\mathbf{x}_1^{[p]}, \dots, \mathbf{x}_k^{[p]}\}$ converge.*

Proof. As we have shown in the proof of Theorem 2.3, the interlacing property in (2.13) implies that

$$\lim_{p \rightarrow \infty} \langle \mathbf{A}(\mathbf{a}^{[p+1]} - \mathbf{a}^{[p]}), (\mathbf{a}^{[p+1]} - \mathbf{a}^{[p]}) \rangle = 0.$$

On one hand, because A is positive semi-definite, we have

$$\lim_{p \rightarrow \infty} \|\mathbf{a}^{[p+1]} - \mathbf{a}^{[p]}\|_F^2 = 2 - 2 \lim_{p \rightarrow \infty} \operatorname{Re}(\Pi_{i=1}^k \langle \mathbf{x}_i^{[p]}, \mathbf{x}_i^{[p+1]} \rangle) = 0.$$

On the other hand, because $|\Pi_{i=1}^k \langle \mathbf{x}_i^{[p]}, \mathbf{x}_i^{[p+1]} \rangle| \leq 1$, $i \in \llbracket k \rrbracket$, it must be that

$$\lim_{p \rightarrow \infty} \langle \mathbf{x}_i^{[p]}, \mathbf{x}_i^{[p+1]} \rangle = 1, \quad i \in \llbracket k \rrbracket.$$

Throughout the algorithm, with respect to arbitrary α , we have required the phase of the first entries of all $\mathbf{x}_{\alpha_1}^{[p]}$ be positive. Therefore, the two unit vectors $\mathbf{x}_i^{[p]}$ and $\mathbf{x}_i^{[p+1]}$, $i \in \llbracket k \rrbracket$, must be gradually aligned as p goes to infinity. In particular,

$$\lim_{p \rightarrow \infty} (\mathbf{x}_i^{[p+1]} - \mathbf{x}_i^{[p]}) = 0, \quad i \in \llbracket k \rrbracket.$$

By the result already established in [31, Lemma 4.10] and [28, Lemma 2.7], the above limiting behavior of increments between two consecutive iterates is sufficient to prove that $\{(\mathbf{x}_1^{[p]}, \dots, \mathbf{x}_k^{[p]})\}$ converges. \square

Algorithm 1 is designed for multipartite rank-1 approximation problem (2.2). Since its convergence theory is complete and its computation is highly efficient, it is tempting to speculate that we can use it as the basic building block in the so-called greedy ALS update scheme for the general problem (1.7). The idea is that, while advancing in $t = 0, 1, \dots$, we repeatedly apply Algorithm 1 to solve a sequence of subproblems of the form

$$(\lambda_j^{[t+1]}, \mathbf{x}_{1,j}^{[t+1]}, \dots, \mathbf{x}_{k,j}^{[t+1]}) := \arg \min_{\substack{\lambda_j \in \mathbb{R}_+, i \in \llbracket k \rrbracket, \\ \mathbf{x}_i \in \mathbb{C}^{m_i}, \|\mathbf{x}_i\|_2=1}} \|\rho_j^{[t+1]} - \lambda_j (\mathbf{x}_1 \mathbf{x}_1^*) \otimes \dots \otimes (\mathbf{x}_k \mathbf{x}_k^*)\|_F^2, \quad j \in \llbracket R \rrbracket, \quad (2.14)$$

where

$$\begin{aligned} \rho_j^{[t+1]} &:= \rho - \sum_{r=1}^{j-1} \lambda_r^{[t+1]} (\mathbf{x}_{1,r}^{[t+1]} \mathbf{x}_{1,r}^{[t+1]*}) \otimes \dots \otimes (\mathbf{x}_{k,r}^{[t+1]} \mathbf{x}_{k,r}^{[t+1]*}) \\ &\quad - \sum_{r=j+1}^R \lambda_r^{[t]} (\mathbf{x}_{1,r}^{[t]} \mathbf{x}_{1,r}^{[t]*}) \otimes \dots \otimes (\mathbf{x}_{k,r}^{[t]} \mathbf{x}_{k,r}^{[t]*}). \end{aligned} \quad (2.15)$$

The matrix $\rho_j^{[t+1]}$ is composed of two parts — the factors $\lambda_r^{[t]}, \mathbf{x}_{1,r}^{[t]}, \dots, \mathbf{x}_{k,r}^{[t]}$, $r \in \llbracket R \rrbracket \setminus \llbracket j \rrbracket$, are available from the t -th iteration and $\lambda_r^{[t+1]}, \mathbf{x}_{1,r}^{[t+1]}, \dots, \mathbf{x}_{k,r}^{[t+1]}$, $r \in \llbracket j-1 \rrbracket$, are newly updated at the $(t+1)$ -th iteration. Nonetheless, such a successive displacement iterative scheme suffers from several issues both computationally and theoretically. First, it is expensive. For each fixed $t = 0, 1, \dots$, we need to sweep through $j \in \llbracket R \rrbracket$, whereas for each fixed $j \in \llbracket R \rrbracket$ we need to solve (2.14) using Algorithm 1 which itself requires iterations. Second, the matrix $\rho_j^{[t+1]}$ is guaranteed only to be Hermitian but is not necessarily positive semi-definite, whereas our convergence analysis for Algorithm 1 relies heavily on the definiteness of the underlying matrix. Third, the foremost challenge in computation is to maintain the critically important probability mixture of separable states in (1.7), i.e., the conditions that

$$\sum_{r=1}^R \lambda_r = 1, \quad \lambda_r \geq 0. \quad (2.16)$$

The $\lambda_j^{[t+1]}$ found by solving the individual problem (2.14), however, does not take this constraint into account as a whole. In fact, since $\rho_j^{[t+1]}$ is a mixture by $\lambda_1^{[t+1]}, \dots, \lambda_{j-1}^{[t+1]}$ and $\lambda_{j+1}^{[t]}, \dots, \lambda_R^{[t]}$ which in general are not even probabilistically related, we have no reason to think that the constraint (2.16) will be satisfied eventually. Enforcing such a condition seems to be a major difficulty in applying the greedy ALS method.

3. Quantum Low-rank Separability Approximation. In contrast to the SVD-based iterative method for the case $R = 1$, in this section we propose a continuous dynamical system approach for the case $R > 1$ when solving the problem (1.7). The dynamical system is based on the complex-valued gradient flow. We observe at least four advantages in such an approach. First, the sum-to-one constraint imposed on the combination coefficients can be built into the dynamical system. Second, any violation of the nonnegativity constraint can easily be detected and fixed. Third, the rank can be automatically adjusted downward and, hence, even if R is wrongly overestimated, it actually helps offer a broader search initially and will be downgraded along the course of integration. Fourth, once the differential equation is in place, the coding is straightforward and any available ODE solver can be used as the numerical integrator.

3.1. Projected Gradient Flow. For convenience, we introduce the abbreviations

$$\begin{cases} \Theta & := \rho - \sum_{r=1}^R \lambda_r (\mathbf{x}_{1,r} \mathbf{x}_{1,r}^* \otimes \cdots \otimes (\mathbf{x}_{k,r} \mathbf{x}_{k,r}^*)) \in \mathbb{C}^{\prod_{i=1}^k m_i \times \prod_{i=1}^k m_i}, \\ \omega_r & := \langle \mathbf{x}_{1,r} \otimes \cdots \otimes \mathbf{x}_{k,r}, \Theta (\mathbf{x}_{1,r} \otimes \cdots \otimes \mathbf{x}_{k,r}) \rangle \in \mathbb{R}, \\ \mathfrak{C}_r & := \text{reshape} (\Theta (\mathbf{x}_{1,r} \otimes \cdots \otimes \mathbf{x}_{k,r}), [m_k, \dots, m_1]) \in \mathbb{C}^{m_k \times \cdots \times m_1}. \end{cases}$$

Note that ω_r and \mathfrak{C}_r vary in $r \in \llbracket R \rrbracket$. Note also that the expressions involve every λ_r , $r \in \llbracket R \rrbracket$. That is, different from the greedy ALS scheme (2.14), we want to adjust the entire array $\{\lambda_1, \dots, \lambda_R\}$ simultaneously. Despite of their seemingly complicated expressions, it will be interesting to find in the following development that Θ , ω_r and \mathfrak{C}_r for the case $R > 1$ generalize the roles of A , λ and \mathfrak{D} discussed in the preceding section for the case $R = 1$, respectively.

Rewrite the objection function in (1.7) as

$$g(\lambda_1, \dots, \lambda_R, \mathbf{x}_{1,1}, \dots, \mathbf{x}_{k,1}, \mathbf{x}_{1,2}, \dots, \mathbf{x}_{k,2}, \dots, \mathbf{x}_{1,R}, \dots, \mathbf{x}_{k,R}) := \langle \Theta, \Theta \rangle = \langle \Theta, \overline{\Theta} \rangle_R. \quad (3.1)$$

It is not difficult to calculate the Wirtinger derivative of the function g with respect to the various variables. We summarize the results as follows:

$$\begin{cases} \frac{\partial g}{\partial \lambda_r} & = -2\omega_r, \\ \frac{\partial g}{\partial \mathbf{x}_{j,r}} & = -2\lambda_r \overline{\mathfrak{C}_r} \otimes_j \left(\bigcirc_{i=k, i \neq j}^1 \mathbf{x}_{i,r} \right), \\ \frac{\partial g}{\partial \overline{\mathbf{x}}_{j,r}} & = -2\lambda_r \mathfrak{C}_r \otimes_j \left(\bigcirc_{i=k, i \neq j}^1 \overline{\mathbf{x}}_{i,r} \right), \end{cases} \quad j \in \llbracket k \rrbracket, r \in \llbracket R \rrbracket. \quad (3.2)$$

Note that the outer product is done specifically in the reverse order. If we denote $\mathbf{x}_{j,r} = \mathbf{u}_{j,r} + \imath \mathbf{v}_{j,r}$ with $\mathbf{u}_{j,r}, \mathbf{v}_{j,r} \in \mathbb{R}^{m_j}$, then by using (2.7) the above Wirtinger gradients (3.2) can be converted to the real gradients as follows:

$$\frac{\partial g}{\partial (\mathbf{u}_{j,r}, \mathbf{v}_{j,r})} = -4\lambda_r \begin{bmatrix} \text{Re}(\mathfrak{C}_r \otimes_j \left(\bigcirc_{i=k, i \neq j}^1 \overline{\mathbf{x}}_{i,r} \right)) \\ \text{Im}(\mathfrak{C}_r \otimes_j \left(\bigcirc_{i=k, i \neq j}^1 \overline{\mathbf{x}}_{i,r} \right)) \end{bmatrix}, \quad j \in \llbracket k \rrbracket, r \in \llbracket R \rrbracket. \quad (3.3)$$

This expression is similar to that in (2.10). Using the same argument as that for deriving (2.12), we arrive at the relationships

$$\omega_r = \mathbf{u}_{j,r}^* \text{Re}(\mathfrak{C}_r \otimes_j \left(\bigcirc_{i=k, i \neq j}^1 \overline{\mathbf{x}}_{i,r} \right)) + \mathbf{v}_{j,r}^* \text{Im}(\mathfrak{C}_r \otimes_j \left(\bigcirc_{i=k, i \neq j}^1 \overline{\mathbf{x}}_{i,r} \right)), \quad r \in \llbracket R \rrbracket. \quad (3.4)$$

Therefore, the projected gradients of objective function g onto the unit sphere S^{2m_j-1} , $j \in \llbracket k \rrbracket$, can be expressed in the condensed form:

$$\text{Proj}_{S^{2m_j-1}} \frac{\partial g}{\partial (\mathbf{u}_{j,r}, \mathbf{v}_{j,r})} = -4\lambda_r (\mathfrak{C}_r \otimes_j \left(\bigcirc_{i=k, i \neq j}^1 \overline{\mathbf{x}}_{i,r} \right) - \omega_r \mathbf{x}_{j,r}), \quad r \in \llbracket R \rrbracket. \quad (3.5)$$

By (3.5), the first-order optimality condition should be that

$$\lambda_r (\mathfrak{C}_r \otimes_j \left(\bigcirc_{i=k, i \neq j}^1 \overline{\mathbf{x}}_{i,r} \right) - \omega_r \mathbf{x}_{j,r}) = 0, \quad j \in \llbracket k \rrbracket, r \in \llbracket R \rrbracket,$$

which resembles that in Lemma 2.1 but is more involved because r varies. By now, we have established a negative gradient flow

$$\begin{cases} \dot{\lambda}_r &= 2\omega_r, \\ \dot{\mathbf{x}}_{j,r} &= 4\lambda_r(\mathfrak{C}_r \otimes_j \left(\bigcirc_{i=1, i \neq j}^k \bar{\mathbf{x}}_{i,r} \right) - \omega_r \mathbf{x}_{j,r}), \end{cases} \quad j \in \llbracket k \rrbracket, r \in \llbracket R \rrbracket, \quad (3.6)$$

whose solution defines a trajectory along which the objective value of (1.7) is gradually decreased. However, thus far there is no guarantee on whether the resulting $\lambda_r(t)$, $r \in \llbracket R \rrbracket$, will satisfy the constraint (2.16). We will modify the differential equation to address this issue while still maintaining the descent property in the next section. We also have to devise an implementation that respects the nonnegativity constraint.

3.2. Modified gradient flow and adaptive strategy. We address the sum-to-one constraint first. Suppose that initially $\lambda_r(0) > 0$, $r \in \llbracket R \rrbracket$, and $\sum_{r=1}^R \lambda_r(0) = 1$. To satisfy the constraint (2.16), it is necessary that

$$\sum_{r=1}^R \dot{\lambda}_r(t) = 0, \quad \text{for all } t \geq 0. \quad (3.7)$$

The dynamical system given in (3.6) alone can hardly meet this condition. We propose to remedy the situation by modifying the flow for $\lambda_r(t)$ to

$$\dot{\lambda}_r = 2(\omega_r - \tilde{\omega}), \quad r \in \llbracket R \rrbracket, \quad (3.8)$$

with $\tilde{\omega} := \frac{\sum_{r=1}^R \omega_r}{R}$, while keeping intact the original governing equations for $\dot{\mathbf{x}}_{j,r}$, $j \in \llbracket k \rrbracket$, $r \in \llbracket R \rrbracket$. By doing it this way, the condition (3.7) is met, but the direction of the flow $\mathbf{x}_{j,r}$, $j \in \llbracket k \rrbracket$, $r \in \llbracket R \rrbracket$ will have been altered. Even so, the following result shows that we still have a descent flow.

LEMMA 3.1. *Let*

$$Z(t) := (\lambda_1(t), \dots, \lambda_R(t), \mathbf{x}_{1,1}(t), \dots, \mathbf{x}_{k,1}(t), \dots, \mathbf{x}_{1,r}(t), \dots, \mathbf{x}_{k,r}(t)) \quad (3.9)$$

denote the flow corresponding to the newly modified differential system described above. Then the objection value of g is descending along the trajectory $Z(t)$.

Proof. We first calculate that

$$\begin{aligned} \frac{dg(Z(t))}{dt} &= \nabla g(Z(t)) \cdot \frac{dZ(t)}{dt} \\ &= \sum_{r=1}^R \frac{\partial g}{\partial \lambda_r} \dot{\lambda}_r + \sum_{j=1}^k \sum_{r=1}^R \left\langle \frac{\partial g}{\partial (\mathbf{u}_{j,r}, \mathbf{v}_{j,r})}, \begin{bmatrix} \dot{\mathbf{u}}_r \\ \dot{\mathbf{v}}_r \end{bmatrix} \right\rangle \\ &= \sum_{r=1}^R \frac{\partial g}{\partial \lambda_r} \dot{\lambda}_r - 16 \sum_{j=1}^k \sum_{r=1}^R \lambda_r^2 \left(\left\| \mathfrak{C}_r \otimes_j \left(\bigcirc_{i=1, i \neq j}^k \bar{\mathbf{x}}_{i,r} \right) \right\|^2 - \omega_r^2 \right). \end{aligned} \quad (3.10)$$

It follows from (3.4) that each term in the last summations is nonnegative. Also,

$$\sum_{r=1}^R \frac{\partial g}{\partial \lambda_r} \dot{\lambda}_r = -4 \sum_{r=1}^R \omega_r (\omega_r - \tilde{\omega}) = -4 \left(\sum_{r=1}^R \omega_r^2 - \frac{1}{R} \left(\sum_{r=1}^R \omega_r \right)^2 \right) \leq 0, \quad (3.11)$$

where the last inequality follows from the Cauchy-Schwarz inequality and the fact that $\omega_r \in \mathbb{R}$, $r \in \llbracket R \rrbracket$. In all, we see that $\frac{dg(Z(t))}{dt} \leq 0$. We mention in passing that the equality in (3.11) holds only if $\omega_r = \frac{1}{\sqrt{R}}$. \square

We next address the task of keeping $\lambda_r \geq 0$, $r \in \llbracket R \rrbracket$. Maintaining nonnegativity in solutions of ordinary differential systems has been widely discussed in the literature. A variety of strategies for enforcing nonnegativity can be found in the literature. See [32] and the references contained therein for a historic review of this subject. For our application, we propose the following mechanism to keep $\lambda_r \geq 0$, $r \in \llbracket R \rrbracket$. The mechanism consists of three components working together:

1. **Event Detection:** By an event we mean that one of these $\lambda_r(t)$, $r \in \llbracket R \rrbracket$, has decreased from a positive value to zero (or near zero) for some t during the integration. It is critical to determine the time \hat{t} when an event occurs up to the prescribed precision. Such a detection machinery can effectively be programmed in any numerical solver. For demonstrate purpose, we shall make use of the existing event function in the Matlab ODE suite to carry out the task.
2. **Rank Reduction:** The event $\lambda_r(\hat{t}) = 0$ indicates two things. First, since $\dot{\lambda}_r(\hat{t}) \leq 0$, any further integration even at a tiny time step is likely to violate the nonnegative constraint. Second, since the term $\lambda_r(\hat{t})(\mathbf{x}_{1,r}(\hat{t})\mathbf{x}_{1,r}^*(\hat{t}) \otimes \cdots \otimes \mathbf{x}_{k,r}(\hat{t})\mathbf{x}_{k,r}^*(\hat{t})) = 0$ is not making any contribution to the objective value of g at the moment, we can drop this term and continue. In doing so, the initial rank R is reduced by one.
3. **Restart:** Once a term is dropped, we use the remaining information $(\lambda_s(\hat{t}), \mathbf{x}_{1,s}(\hat{t}), \dots, \mathbf{x}_{k,s}(\hat{t}))$, $s \in \llbracket R \rrbracket \setminus \{r\}$, as the initial value to restart the integration. In this way, the objective value is ratcheted at the current value and can only continue to go down after the restart.

Recall that estimating a proper R is always difficult in low-rank approximation. Starting with a larger rank R might seem redundant and wasteful initially, but it provides the flexibility of searching multiple directions for a better solution. The mechanism described above serves as a means to filter out unneeded factors.

3.3. Convergence. The limiting behavior of a gradient dynamics is well studied in the literature. In particular, counterexamples have been found to evince that not all gradient flow will converge. For completion, we now argue that our gradient flow, even with the modification (3.8), will converge to a singleton point.

If we separate each $\mathbf{x}_{j,r}$ into real and imaginary parts, the right-hand side of our differential system can be regarded as a polynomial system in a total of $(2 \sum_{i=1}^k m_i + 1)R$ real variables. Without loss of the original sense, let ξ denote the vector of all real variables and abridge the differential system as a negative gradient flow

$$\frac{d\xi}{dt} = -\nabla F(\xi) \quad (3.12)$$

for some abstract objective function $F(\xi)$ in ξ . Being polynomials in ξ , the vector field $\nabla F(\xi)$ is real analytic in ξ . By construction, $\xi(t)$ is also bounded. It follows that the set of accumulation points

$$\omega(\xi(0)) := \left\{ \tilde{\xi} \in \mathbb{R}^n \mid \mathbf{x}(t_\nu) \rightarrow \tilde{\xi} \text{ for some sequence } t_\nu \rightarrow \infty \right\} \quad (3.13)$$

is a non-empty, compact, and connected subset of stationary points, i.e., $\nabla F(\tilde{\xi}) = 0$.

Recall the following Łojasiewicz gradient inequality [33, 34].

THEOREM 3.2. *Suppose that $F : U \rightarrow \mathbb{R}$ is a real analytic function in an open set $U \subset \mathbb{R}^n$ and $\tilde{\xi} \in U$. Then there exist a neighborhood W of $\tilde{\xi}$, constants $c > 0$ and $\theta \in [0, 1)$, such that the inequality*

$$\|\nabla F(\xi)\| \geq c \|F(\xi) - F(\tilde{\xi})\|^\theta$$

holds for all $\xi \in W$.

An important consequence of the Łojasiewicz gradient inequality is that the trajectory of an analytic gradient flow is necessary of finite length. The following result is readily applicable to our differential system and implies that the flow $\{\lambda_r(t), \mathbf{x}_{j,r}(t)\}$, $j \in \llbracket k \rrbracket$, $r \in \llbracket R \rrbracket$ converge to a singleton stationary point. See [35, Theorem 2.2] and the lecture note [36] for its proof.

THEOREM 3.3. *Suppose that $F : U \rightarrow \mathbb{R}$ is real analytic in an open set $U \subset \mathbb{R}^n$. Then for any bounded semi-orbit of (3.12)*

$$\xi(t) \rightarrow \tilde{\xi} \quad \text{as } t \rightarrow \infty$$

for some $\tilde{\xi} \in U$.

4. Numerical Experiment. Quantifying the entanglement of a mixed state by finding its nearest separable state is a challenging problem. Even the problem of determining whether a given state is separable or not is NP-hard [37, 38]. However, we must not misconceive that an NP-hard problem is forever hopeless and untouchable. A practical example is the problem that expressing $\sqrt{2}$ is NP-hard in theory because it requires infinite complexity on a Turing machine but we do have polynomial-time algorithms to approximate it to any finite precision. Likewise, the NP-hardness encountered in entanglement quantification does not imply that we cannot approximately solve the problem by numerical means. Thus far, we have described an SVD-based iterative method for the rank-1 k -partite problem (2.2) and a gradient flow approach for the more complicated rank- R k -partite problem. In this section we carry out some numerical experiments to test the effectiveness of our methods.

Example 1. In the first part of this experiment, we produce the target matrix

$$\Xi_1 = (\boldsymbol{\eta}_1 \otimes \boldsymbol{\eta}_2 \otimes \boldsymbol{\eta}_3 \otimes \boldsymbol{\eta}_4)(\boldsymbol{\eta}_1 \otimes \boldsymbol{\eta}_2 \otimes \boldsymbol{\eta}_3 \otimes \boldsymbol{\eta}_4)^*$$

with randomly generated unit vectors $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_4 \in \mathbb{C}^5$. Therefore, the target matrix $\rho \in \mathbb{C}^{625 \times 625}$ is already separable and is of rank one. We test Algorithm 1 with randomly selected $\boldsymbol{\alpha}$ as well as its analogue where $\boldsymbol{\alpha}$ is varied cyclically. The iteration is terminated whenever the stopping criterion

$$\left\| \left[\mathfrak{D}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_k) \otimes_j \bigcirc_{i=1, i \neq j}^4 \bar{\mathbf{x}}_i - \lambda \mathbf{x}_j \right]_{j=1, \dots, 4} \right\|_F < 10^{-10}. \quad (4.1)$$

is met. We repeat our experiments 20 times with distinct randomly generated starting points. As a nonlinear optimization problem, the limit points depend on the starting points and may differ from the original generators $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_4 \in \mathbb{C}^5$. It should be more feasible if we gauge the quality of the approximation not by a comparison with the original generators but by the product

$$\hat{\rho} = (\hat{\mathbf{x}}_1 \otimes \hat{\mathbf{x}}_2 \otimes \hat{\mathbf{x}}_3 \otimes \hat{\mathbf{x}}_4)(\hat{\mathbf{x}}_1 \otimes \hat{\mathbf{x}}_2 \otimes \hat{\mathbf{x}}_3 \otimes \hat{\mathbf{x}}_4)^*$$

based on the limit point $(\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \hat{\mathbf{x}}_3, \hat{\mathbf{x}}_4)$ of the iterates. We measure the quantity

$$\text{Error} := \|\Xi_1 - \hat{\rho}\|_F.$$

The test results in terms of the averages of errors, numbers of iterations, and CPU time in seconds of the 20 runs are tabulated in Table 4.1. Since Algorithm 1 utilizes only the first-order derivative information, its rate of convergence should be at most linear. Nonetheless, for this problem of decomposing a 625×625 density matrix as the tensor product of four pure state density matrices in \mathbb{C}^5 , the empirical data seem to suggest that our SVD-based approach can be effective in precision and efficient in time when calculating the optima.

$\boldsymbol{\alpha}$	Error	Iteration	MinTime	MaxTime	AveTime
cyclic	9.6389×10^{-16}	4	4.0921×10^{-3}	1.1574×10^{-1}	1.0126×10^{-2}
random	3.5439×10^{-15}	4.1	2.9594×10^{-3}	2.4689×10^{-2}	5.3788×10^{-3}

Table 4.1: Average errors, numbers of iterations, and CPU time in seconds on 20 runs by Algorithm 1 for ρ .

In practice, the exact rank of a given entangled state is not known. Indeed, for almost all information gathering devices, it is inevitable that the data collected contain noise. The presence of even a small amount of noise to a low-rank matrix will break up the low rank. In the second part of this experiment, we mete out the perturbation in a controlled way and calculate the rank-1 approximation of ρ under noise. Specifically, we use the perturbed matrix

$$\rho_\sigma = \Xi_1 + \sigma(\gamma - \Xi_1) \quad (4.2)$$

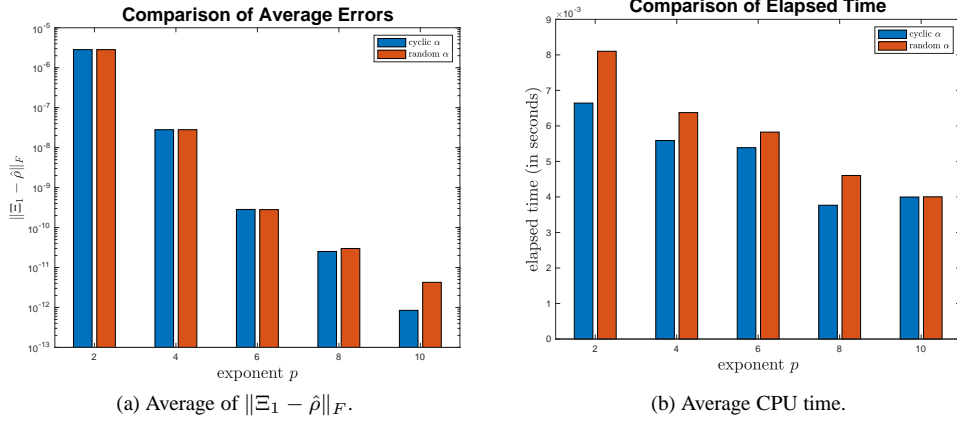


Figure 4.1: Average errors and CPU time on 20 random runs by Algorithm 1 for ρ_σ .

as the target matrix, where γ is a randomly generated but fixed density matrix and $\sigma = 10^{-p}$, $p = 2, 4, 6, 8, 10$, represents the intensity of the noise. In this particular experiment, we already know $\|\gamma - \Xi_1\| \approx 1.0001$, so $\sigma \approx \|\rho_\sigma - \Xi_1\|$. Observe that ρ_σ is still a density matrix and is of full rank in general, but ρ_σ is made more of the rank-1 matrix Ξ_1 than of the full rank matrix γ . For this reason, we find in Figure 4.1(a) that the ultimate rank-1 approximation $\hat{\rho}$ for ρ_σ is in fact closer to the rank-1 matrix Ξ_1 than to the full rank matrix ρ_σ . We also find that the smaller the perturbation σ , the closer $\hat{\rho}$ is to Ξ_1 . The CPU time is little affected.

Example 2. Despite the success of the SVD-based iterative method for the case $R = 1$, it is difficult to generalize to the greedy method for the case $R > 1$. In the first part of this second experiment, we demonstrate that the differential system approach can identify a proper rank by the mechanism described in Section 3.2. Consider a separable rank-2 target matrix:

$$\Xi_2 = \sum_{r=1}^2 \lambda_r (\boldsymbol{\eta}_{1,r} \boldsymbol{\eta}_{1,r}^*) \otimes (\boldsymbol{\eta}_{2,r} \boldsymbol{\eta}_{2,r}^*) \otimes (\boldsymbol{\eta}_{3,r} \boldsymbol{\eta}_{3,r}^*) \otimes (\boldsymbol{\eta}_{4,r} \boldsymbol{\eta}_{4,r}^*), \quad (4.3)$$

where $\boldsymbol{\eta}_{i,r} \in \mathbb{C}^5$, $i = 1, \dots, 4$, $r = 1, 2$, are randomly generated unit vectors and $\lambda_r > 0$, $r = 1, 2$, satisfies $\sum_{r=1}^2 \lambda_r = 1$. Pretending that we do not know of the exact low rank of Ξ_2 initially, we start off with $R = 4$ with the hope the exact rank of Ξ_2 will be found eventually.

Utilizing the existing routine `ode15s` in `Matlab` as the integrator, we turn on the option `event` and set the local error tolerance at `AbSTol` = 10^{-10} and `RelTol` = 10^{-10} . We follow four trajectories, each with a different set of starting points. The evolution of the objective values when following these trajectories is plotted in Figure 4.2(a). As can be seen, though the integral curves follow different trajectories and might end up with distinct stationary points, the ultimate objective values can be regarded as nearly zero within the prescribed tolerance. Evidence about the preservation of the sum-to-one property, even with the overestimated R , is plotted in Figure 4.2(b). Though their markings might be smeared due to the proximity of the curves, the red circles in both graphs in Figure 4.2 indicate that an event has been detected and, hence, the value R is reduced by one. There are two red circles in each curve, so we know that the original low rank has been found.

Similar to the experiment done in Example 1, we next investigate how the noise affects our dynamical approach's performance. We deal out the perturbation to Ξ_2 in exactly the same way as in (4.2) to produce the target matrix ρ_σ . Let $\hat{\rho}(t)$ denote the numerical solution to our differential equation. The evolution of $\|\Xi_2 - \hat{\rho}(t)\|_F$ and $\|\rho_\sigma - \hat{\rho}(t)\|_F$ in response to different levels of perturbation strength σ are plotted in Figure 4.3. In contrast to the case of $R = 1$ in Example 1, we have observed that the rank-2 separable state Ξ_2 is more sensitive to perturbation in the sense that $\|\Xi_2 - \hat{\rho}_\sigma(t)\|_F \approx \|\rho_\sigma - \hat{\rho}_\sigma(t)\|_F$ if $\sigma \geq 10^{-6}$. That is, the flow $\hat{\rho}_\sigma(t)$ is about equal distance to both the rank-2 state Ξ_2 and the full rank state ρ_σ if the perturbation σ is lightly too large.

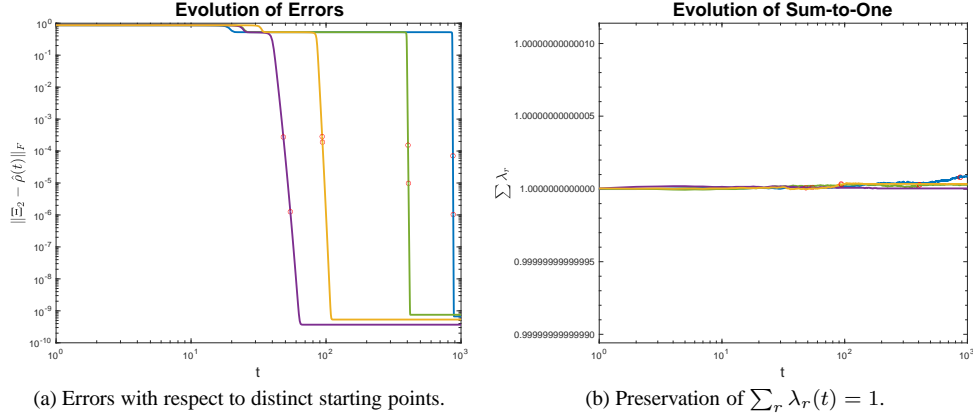


Figure 4.2: Convergence, event detection, and sum-to-one constraint for approximating exact Ξ_2 .

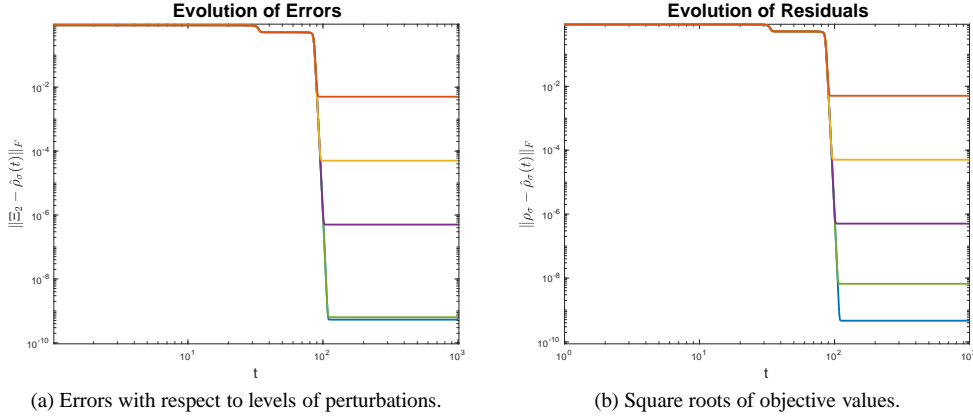


Figure 4.3: Rank-4 approximation for perturbed matrices ρ_σ of Ξ_2 , $\sigma = 10^{-p}$, $p = 2, 4, 6, 8, 10$.

Different from the case $R = 1$, two observations are worth noting. First, $\hat{\rho}(t)$ does not show a preference to Ξ_2 . Second, we start with $R = 4$, but not all trajectories encounter an event to reduce the rank. The latter phenomenon leads us to reconsider the definition of the event,

Instead of defining an event rigorously only at $\lambda_r(\hat{t}) = 0$, we may declare that an event happens whenever $|\lambda_r(\hat{t})| < \epsilon$ at a preselected ϵ . The rationale is that since $\|\mathbf{x}_{j,r}(t)\|_2 = 1$, $j \in \llbracket k \rrbracket$, for all t , the contribution of the term $\lambda_r(\hat{t})(\mathbf{x}_{1,r}(\hat{t})\mathbf{x}_{1,r}(\hat{t})^* \otimes \cdots \otimes (\mathbf{x}_{k,r}(\hat{t})\mathbf{x}_{k,r}(\hat{t})^*))$ to the overall summation is at most $|\lambda_r(\hat{t})|$. If the current approximation is several orders larger than ϵ , then maybe it is plausible to ignore the term whose contribution is at most ϵ . We experiment with this more relaxed rank reduction mechanism on the perturbed matrix ρ_σ with $\sigma = 10^{-4}$ by choosing $\epsilon = .5 * 10^{-4}$. We start with $\lambda_r(0) = \frac{1}{R}$, $r \in \llbracket R \rrbracket$, to give every component an equal chance to compete for survival. Since now it is easier to qualify an event, we set $R = 40$ with the hope of broadening the search in many more directions. The red line at the top of Figure 4.4(a) shows the preservation of the sum-to-one property. We also see that at $t \approx 100$, many $\lambda_r(t)$'s begin to diminish and eventually trigger the event mechanism. The rank reduction from $R = 40$ to $R = 2$ happens quickly within a small window of t , as can be seen from the cluster of red circles in Figure 4.4(b). While the rank is being reduced, we see that the objective values continue to decrease until a local minimum is found. This experiment supports the mechanism of relaxed event qualification, but in practice we usually do not have a priori knowledge of the proper extent of relaxation.

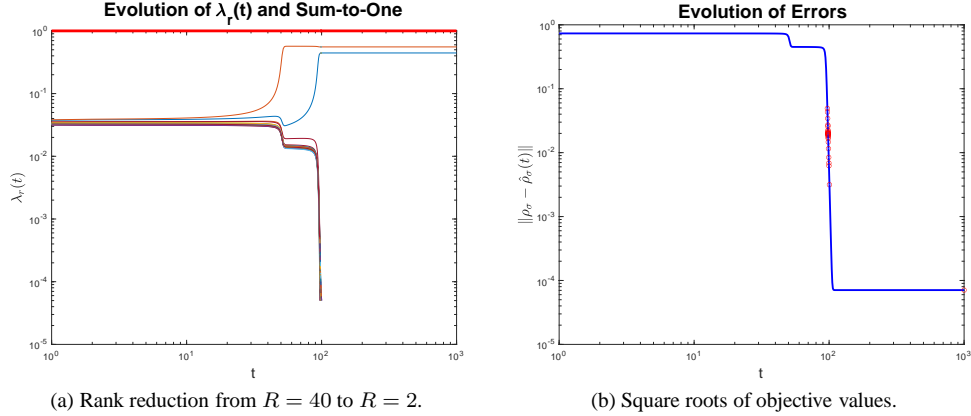


Figure 4.4: Rank reduction for perturbed matrix ρ_σ with $\sigma = 10^{-4}$ and relaxed event qualification.

Example 3. Let $\Xi_3 \in \mathbb{C}^{32 \times 32}$ be a randomly generated positive definite matrix. Suppose that Ξ_3 is regarded as a density matrix in a 5-qubit system. Recall that each qubit counts as an element in \mathbb{C}^2 . Let the notation $\llbracket p_1, p_2, \dots, p_\ell \rrbracket$ denote the composite system $(\mathbb{C}^2)^{\otimes p_1} \otimes (\mathbb{C}^2)^{\otimes p_2} \otimes \dots \otimes (\mathbb{C}^2)^{\otimes p_\ell}$. There are seven ways to split the number 5 as a sum of nonnegative integers. (The number of partitions for a d -qubit is the Sloane sequence A000041.) Thus we may consider partial separability approximations associated with the group assignments: $\llbracket 1, 1, 1, 1, 1 \rrbracket$, $\llbracket 2, 1, 1, 1 \rrbracket$, $\llbracket 2, 2, 1 \rrbracket$, $\llbracket 3, 1, 1 \rrbracket$, $\llbracket 3, 2 \rrbracket$, $\llbracket 4, 1 \rrbracket$, and $\llbracket 5 \rrbracket$.

The case $\llbracket 5 \rrbracket$ amounts to the rank- R approximation in the conventional sense of linear algebra. The problem can be resolved directly by the singular value decomposition of Ξ_3 . It is a well-known fact from the Eckart-Young-Mirsky theorem that, for each $R \leq 32$, the truncated singular value decomposition gives rise to the globally best rank- R approximation to Ξ_3 [39, 19].

Starting from $R = 50$, we apply the gradient flow approach to approximate Ξ_3 over the other six types of multipartite systems specified by the group assignments. Figure 4.5 shows the evolution of errors $\|\Xi_3 - \hat{\rho}(t)\|_F$. We see that as the number of splits decreases, the errors are reduced correspondingly. This observation is expected because, for example, the case $\llbracket 3, 1, 1 \rrbracket$ can be considered as a more restrictive structure of $\llbracket 3, 2 \rrbracket$ and, hence, it should have higher errors. In the order from $\llbracket 1, 1, 1, 1, 1 \rrbracket$ to $\llbracket 4, 1 \rrbracket$, we find that the final reduced ranks are 41, 46, 41, 43, 41, and 36, respectively. The fluctuation of the final ranks might indicate that we have found a local solution only, but the general trend is that the higher the separability is involved, the more demanding is the computation. It is interesting to note that the singular value decomposition of a generic Ξ_3 requires exactly $R = 32$ for a complete decomposition of Ξ_3 . The fact that our gradient approach for the $\llbracket 4, 1 \rrbracket$ -type approximation reduces the final rank to 36 seems to evince that the rank reduction mechanism works reasonably well.

Example 4. In this experiment, we apply our gradient flow to a realistic problem. We briefly describe some background information before carrying out the experiment. In the quantum information theory, the so-called Greenberger-Horne-Zeilinger state (GHZ state) [40, 41]

$$|GHZ\rangle = \frac{1}{\sqrt{2}}(|0\rangle^{\otimes k} + |1\rangle^{\otimes k})$$

is a quantum state that involves the entanglement of at least three subsystems, i.e., $k \geq 3$. Because they exhibit some extremely non-classical properties, GHZ states are used in several protocols in quantum communication and cryptography. In this example, we consider the simplest case $k = 3$ and a density matrix W_σ of the form

$$W_\sigma := (1 - \sigma) |GHZ\rangle \langle GHZ| + \sigma \frac{1}{8} I_8, \quad 0 \leq \sigma \leq 1,$$

which represents a probabilistic mixture of $|GHZ\rangle$ with the operator $\frac{1}{8} I_8$. The mixed state W_σ is known as the

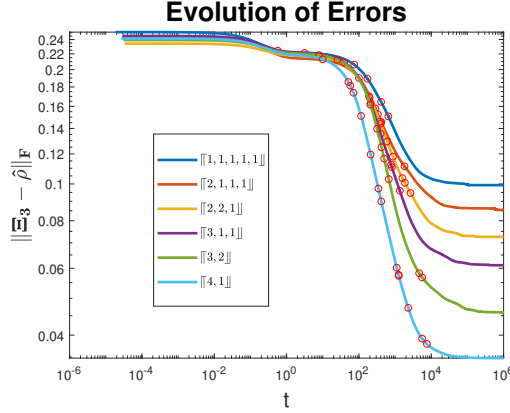


Figure 4.5: Low rank approximation of a 5-qubit system.

generalized Werner state which has found applications in the robustness of entanglement [42], NMR quantum computation [43], and purification schemes for entangled states [44].

The matrix representation of W_σ can be expressed as

$$W_\sigma = (1 - \sigma) \begin{bmatrix} \frac{1}{2} & 0 & \cdots & 0 & \frac{1}{2} \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & 0 \\ \frac{1}{2} & 0 & \cdots & 0 & \frac{1}{2} \end{bmatrix} + \sigma \frac{1}{8} \begin{bmatrix} 1 & 0 & \cdots & \cdots & 0 \\ 0 & 1 & 0 & \cdots & \vdots \\ \vdots & 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & \cdots & \cdots & 0 & 1 \end{bmatrix} \in \mathbb{R}^{8 \times 8}.$$

It has been shown in theory that W_σ is (totally) separable if and only if $\frac{4}{5} \leq \sigma \leq 1$ [45]. Therefore, by adjusting σ we have a test case to explore the rank reduction mechanism and to estimate a possibly optimal rank.

Since the exact low rank of W_σ is not known a priori, we start out our experiments from $R = 15$ for the three choices of $\sigma = 2/3, 4/5$ and $6/7$. We are interested in observing two phenomena in each choice. For the case $\sigma = 6/7$, how low can the rank be reduced and can the separability be achieved? What will happen to the case $\sigma = 2/3$ which is not separable? Can the rank be reduced at all while the objective value is decreased? The borderline case $\sigma = 4/5$ is most curious. It is separable in theory, but will its rank be the same as that for $\sigma = 6/7$ at total separation?

Starting from the same randomly generated unit vectors $\mathbf{x}_{1,r}, \mathbf{x}_{2,r}, \mathbf{x}_{3,r} \in \mathbb{C}^2$, $r \in \llbracket 15 \rrbracket$, and using the rigorous event qualification, we plot the evolution trajectories of the residuals together with red circles whenever an event has been detected in Figure 4.6. The monotone decreasing property guaranteed by our theory is clearly manifested in these curves. The answers to the above questions for the cases $\sigma = 2/3$ and $\sigma = 6/7$ are also clear. It is estimated that the flow for $\sigma = 2/3$ reaches its local solution much sooner (at $t \approx 300$) than the flow for $\sigma = 6/7$ reaches its total separability (at $t \approx 34000$). These limiting behaviors strongly support that W_σ is entangled if $\sigma = 2/3$ and is separable if $\sigma = 6/7$. For the case $\sigma = 4/5$, we estimate that the residual decreases at approximately the rate $O(1/t)$. That is, while the residuals keep going down, it converges at a very slow pace. The behavior suggests that the case $\sigma = 4/5$ is separable, but it is much harder to find its components.

Note that in the extreme case $\sigma = 1$,

$$W_1 = \frac{1}{8} I_8 = \frac{1}{8} \sum_{i=0}^7 |i\rangle \langle i|,$$

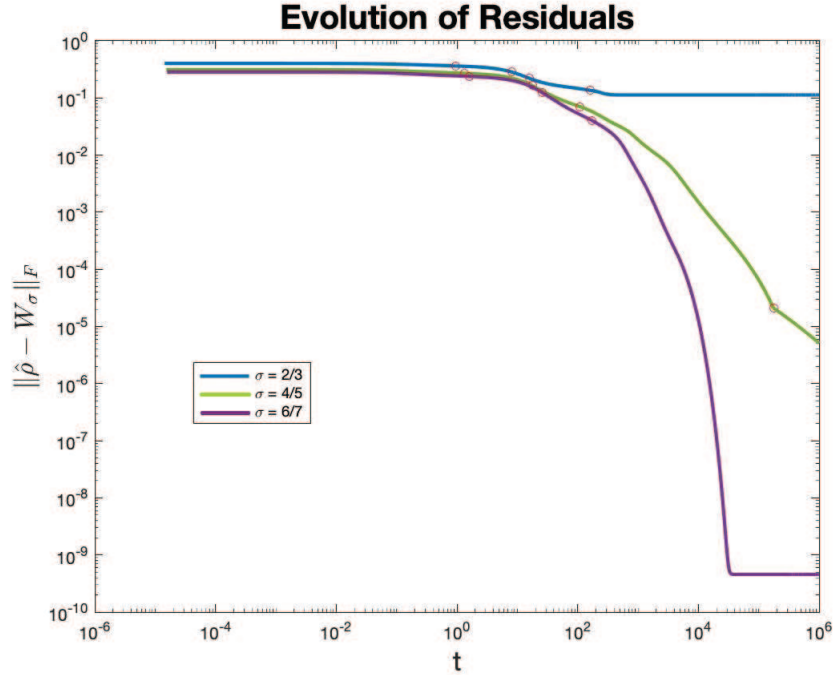


Figure 4.6: Low rank approximation of entangled density matrices.

where $|i\rangle$ denotes the 3-qubit whose binary representation is equal to i . That is, the optimal rank for W_1 is $R = 8$. Counting the events for cases $\sigma = 2/3, 4/5$, and $6/7$, we observe that the reduced ranks are 11, 11, and 12, respectively. While our rank reduction mechanism is tested to work for arbitrary σ , it remains an open question on whether these are the optimal ranks. It is also unclear whether the trajectory for the case $\sigma = 4/5$ will have additional events in later stage of integration.

5. Conclusion. Quantum technologies have been rapidly advanced with the urgent need to create more complex and powerful quantum computers. The technologies, if fully developed, will have far-reaching applications including, for example, as critical as superior analytics capabilities or as practical as better battery life. At the crux of quantum computing is the understanding and control of quantum entanglement, which has already attracted many research endeavors. This paper is concerned with computing numerically the low rank separable approximation of a given entangled multipartite system, which might be used as a computational tool for gauging the quality of entanglement of quantum states.

The notion of quantum mechanics is generally described in physics terms, but there is rich mathematics involved. This work employs a synthesis of techniques from linear algebra, optimization, and dynamical system to tackle the entanglement certification problem numerically. All discussions are over the complex field, so the methods are readily transferable to real-world problems. For rank-1 approximations, the SVD-based iterative method is shown to be efficient and effective. For higher rank approximations, this work derives a complex-valued differential system that not only guarantees global convergence but also is capable of maintaining a probabilistic ensemble of pure states while dynamically estimating a proper rank in the ensemble.

REFERENCES

- [1] A. EINSTEIN, B. PODOLSKY, AND N. ROSEN, *Can quantum-mechanical description of physical reality be considered complete?*, Phys. Rev., 47 (1935), pp. 777–780, <https://doi.org/10.1103/PhysRev.47.777>.

- [2] N. FRIIS, G. VITAGLIANO, M. MALIK, AND M. HUBER, *Entanglement certification from theory to experiment*, Nature Reviews Physics, 1 (2019), pp. 72–87, <https://doi.org/10.1038/s42254-018-0003-5>.
- [3] O. GÜHNE AND G. TÓTH, *Entanglement detection*, Phys. Rep., 474 (2009), pp. 1–75, <https://doi.org/10.1016/j.physrep.2009.02.004>.
- [4] R. HORODECKI, P. HORODECKI, M. HORODECKI, AND K. HORODECKI, *Quantum entanglement*, Rev. Mod. Phys., 81 (2009), pp. 865–942, <https://doi.org/10.1103/RevModPhys.81.865>.
- [5] G. DAHL, J. M. LEINAAS, J. MYRHEIM, AND E. OVRUM, *A tensor product matrix approximation problem in quantum physics*, Linear Algebra Appl., 420 (2007), pp. 711–725, <https://doi.org/10.1016/j.laa.2006.08.026>.
- [6] S.-H. KYE, *Necessary conditions for optimality of decomposable entanglement witnesses*, Rep. Math. Phys., 69 (2012), pp. 419–426, [https://doi.org/10.1016/S0034-4877\(13\)60007-5](https://doi.org/10.1016/S0034-4877(13)60007-5).
- [7] W. THIRRING, R. A. BERTLMANN, P. KÖHLER, AND H. NARNHOFER, *Entanglement or separability: the choice of how to factorize the algebra of a density matrix*, The European Physical Journal D, 64 (2011), pp. 181–196, <https://doi.org/10.1140/epjd/e2011-20452-1>.
- [8] S. AARONSON, *Quantum computing since Democritus*, Cambridge University Press, Cambridge, 2013, <https://doi.org/10.1017/CBO9780511979309>.
- [9] F. HIAI AND D. PETZ, *Introduction to matrix analysis and applications*, Universitext, Springer, Cham; Hindustan Book Agency, New Delhi, 2014, <https://doi.org/10.1007/978-3-319-04150-6>.
- [10] M. NAKAHARA AND T. OHMI, *Quantum computing: From linear algebra to physical realizations*, CRC Press, Boca Raton, FL, 2008, <https://doi.org/10.1201/9781420012293>.
- [11] M. A. NIELSEN AND I. L. CHUANG, *Quantum Computation and Quantum Information*, Cambridge University Press, 2010, <https://doi.org/10.1017/CBO9780511976667>.
- [12] R. KARAM, *Why are complex numbers needed in quantum mechanics? some answers for the introductory level*, American Journal of Physics, 88 (2020), pp. 39–45, <https://doi.org/10.1119/10.0000258>.
- [13] M.-O. RENOUE, D. TRILLO, M. WEILENMANN, T. P. LE, A. TAVAKOLI, N. GISIN, A. ACÍN, AND M. NAVASCUÉS, *Quantum theory based on real numbers can be experimentally falsified*, Nature, (2021), pp. 1–5, <https://doi.org/10.1038/s41586-021-04160-4>.
- [14] K. CHEN AND L.-A. WU, *A matrix realignment method for recognizing entanglement*, Quantum Inf. Comput., 3 (2003), pp. 193–202, <https://doi.org/10.26421/QIC3.3-1>.
- [15] R. F. WERNER, *Quantum states with einstein-podolsky-rosen correlations admitting a hidden-variable model*, Phys. Rev. A, 40 (1989), pp. 4277–4281, <https://doi.org/10.1103/PhysRevA.40.4277>.
- [16] M. HORODECKI, P. HORODECKI, AND R. HORODECKI, *Mixed-State Entanglement and Quantum Communication*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2001, pp. 151–195, https://doi.org/10.1007/3-540-44678-8_5.
- [17] L. CHEN, M. AULBACH, AND M. HAJDUŠEK, *Comparison of different definitions of the geometric measure of entanglement*, Phys. Rev. A, 89 (2014), p. 042305, <https://doi.org/10.1103/PhysRevA.89.042305>.
- [18] J. M. LEINAAS, J. MYRHEIM, AND E. OVRUM, *Geometrical aspects of entanglement*, Phys. Rev. A (3), 74 (2006), pp. 012313, 13, <https://doi.org/10.1103/PhysRevA.74.012313>.
- [19] G. H. GOLUB AND C. F. VAN LOAN, *Matrix computations*, Johns Hopkins Studies in the Mathematical Sciences, Johns Hopkins University Press, Baltimore, MD, fourth ed., 2013.
- [20] R. WEBSTER, *Convexity*, Oxford Science Publications, The Clarendon Press, Oxford University Press, New York, 1994.
- [21] Z.-A. JIA, R. ZHAI, S. YU, Y.-C. WU, AND G.-C. GUO, *Hierarchy of genuine multipartite quantum correlations*, Quantum Inf. Process., 19 (2020), pp. Paper No. 419, 13, <https://doi.org/10.1007/s11128-020-02922-z>.
- [22] T. G. KOLDA AND B. W. BADER, *Tensor decompositions and applications*, SIAM Rev., 51 (2009), pp. 455–500, <https://doi.org/10.1137/07070111X>, <http://dx.doi.org/10.1137/07070111X>.
- [23] N. VERVLIET, O. DEBALS, L. SORBER, M. VAN BAREL, AND L. DE LATHAUWER, *Tensorlab 3.0*, Mar. 2016, <https://www.tensorlab.net>. Available online.
- [24] M. T. CHU AND M. M. LIN, *Nonlinear power-like and SVD-like iterative schemes with applications to entangled bipartite rank-1 approximation*, SIAM J. Sci. Comput., 43 (2021), pp. S448–S474, <https://doi.org/10.1137/20M1336059>.
- [25] C. F. VAN LOAN, *Structured matrix problems from tensors*, in Exploiting hidden structure in matrix computations: algorithms and applications, vol. 2173 of Lecture Notes in Math., Springer, Cham, 2016, pp. 1–63.
- [26] Y. GUAN, M. T. CHU, AND D. CHU, *SVD-based algorithms for the best rank-1 approximation of a symmetric tensor*, SIAM J. Matrix Anal. Appl., 39 (2018), pp. 1095–1115, <https://doi.org/10.1137/17M1136699>.
- [27] W. WIRTINGER, *Zur formalen Theorie der Funktionen von mehr komplexen Veränderlichen*, Math. Ann., 97 (1927), pp. 357–375, <https://doi.org/10.1007/BF01447872>.
- [28] Y. GUAN, M. T. CHU, AND D. CHU, *Convergence analysis of an SVD-based algorithm for the best rank-1 tensor approximation*, Linear Algebra Appl., 555 (2018), pp. 53–69, <https://doi.org/10.1016/j.laa.2018.06.006>.
- [29] C. B. GARCÍA AND T.-Y. LI, *On the number of solutions to polynomial systems of equations*, SIAM J. Numer. Anal., 17 (1980), pp. 540–546, <https://doi.org/10.1137/0717046>.
- [30] A. J. SOMMESE AND C. W. WAMPLER, *The Numerical Solution of Systems of Polynomials Arising in Engineering and Science*, WORLD SCIENTIFIC, 2005, <https://doi.org/10.1142/5763>.
- [31] J. J. MORÉ AND D. C. SORENSEN, *Computing a trust region step*, SIAM J. Sci. Statist. Comput., 4 (1983), pp. 553–572, <https://doi.org/10.1137/0904038>.
- [32] L. F. SHAMPINE, S. THOMPSON, J. A. KIERZENKA, AND G. D. BYRNE, *Non-negative solutions of ODEs*, Appl. Math. Comput., 170 (2005), pp. 556–569, <https://doi.org/10.1016/j.amc.2004.12.011>.
- [33] R. CHILL, *On the lojasiewicz-Simon gradient inequality*, J. Funct. Anal., 201 (2003), pp. 572–601, [https://doi.org/10.1016/S0022-247X\(03\)00000-0](https://doi.org/10.1016/S0022-247X(03)00000-0).

1016/S0022-1236(02)00102-7.

- [34] S. ŁOJASIEWICZ, *Une propriété topologique des sous-ensembles analytiques réels*, in *Les Équations aux Dérivées Partielles* (Paris, 1962), Éditions du Centre National de la Recherche Scientifique, Paris, 1963, pp. 87–89.
- [35] P.-A. ABSIL, R. MAHONY, AND B. ANDREWS, *Convergence of the iterates of descent methods for analytic cost functions*, *SIAM J. Optim.*, 16 (2005), pp. 531–547, <https://doi.org/10.1137/040605266>.
- [36] M. PIERRE, *Quelques applications de l'inégalité de Lojasiewicz à des discrétisations d'EDP*. SMAI, 2011. <http://smai.emath.fr/smai2011/slides/mpierre/Slides.pdf>.
- [37] S. GHARIBIAN, *Strong NP-hardness of the quantum separability problem*, *Quantum Information & Computation*, 10 (2010), pp. 343–360.
- [38] L. GURVITS, *Classical complexity and quantum entanglement*, *J. Comput. System Sci.*, 69 (2004), pp. 448–484, <https://doi.org/10.1016/j.jcss.2004.06.003>.
- [39] A. EKERT AND P. L. KNIGHT, *Entangled quantum systems and the schmidt decomposition*, *American Journal of Physics*, 63 (1995), pp. 415–423, <https://doi.org/10.1119/1.17904>.
- [40] C. ELTSCHKA AND J. SIEWERT, *Entanglement of three-qubit greenberger-horne-zeilinger-symmetric states*, *Phys. Rev. Lett.*, 108 (2012), p. 020502, <https://doi.org/10.1103/PhysRevLett.108.020502>.
- [41] D. M. GREENBERGER, M. A. HORNE, AND A. ZEILINGER, *Going Beyond Bell's Theorem*, Springer Netherlands, Dordrecht, 1989, pp. 69–72, https://doi.org/10.1007/978-94-017-0849-4_10.
- [42] G. VIDAL AND R. TARRACH, *Robustness of entanglement*, *Phys. Rev. A*, 59 (1999), pp. 141–155, <https://doi.org/10.1103/PhysRevA.59.141>.
- [43] S. L. BRAUNSTEIN, C. M. CAVES, R. JOZSA, N. LINDEN, S. POPESCU, AND R. SCHACK, *Separability of very noisy mixed states and implications for nmr quantum computing*, *Phys. Rev. Lett.*, 83 (1999), pp. 1054–1057, <https://doi.org/10.1103/PhysRevLett.83.1054>.
- [44] M. MURAO, M. B. PLENIO, S. POPESCU, V. VEDRAL, AND P. L. KNIGHT, *Multiparticle entanglement purification protocols*, *Phys. Rev. A*, 57 (1998), pp. R4075–R4078, <https://doi.org/10.1103/PhysRevA.57.R4075>.
- [45] W. DÜR AND J. I. CIRAC, *Classification of multiqubit mixed states: Separability and distillability properties*, *Phys. Rev. A*, 61 (2000), p. 042314, <https://doi.org/10.1103/PhysRevA.61.042314>.