

# OPTIMAL $\ell^1$ RANK ONE MATRIX DECOMPOSITION

RADU BALAN, KASSO A. OKOUDJOU, MICHAEL RAWSON, YANG WANG, AND RUI ZHANG

ABSTRACT. In this paper we consider the decomposition of positive semidefinite matrices as a sum of rank one matrices. We introduce and investigate the properties of various measures of optimality of such decompositions. For some classes of positive semidefinite matrices we give explicitly these optimal decompositions. These classes include diagonally dominant matrices and certain of their generalizations,  $2 \times 2$ , and a class of  $3 \times 3$  matrices.

## 1. INTRODUCTION

The finite dimensional matrix factorization problem that we shall investigate was partially motivated by a related infinite dimensional problem, which we briefly recall.

Suppose that  $\mathbb{H}$  is an infinite-dimensional separable Hilbert space, with norm  $\|\cdot\|$  and inner product  $\langle \cdot, \cdot \rangle$ . Let  $\mathcal{I}_1 \subset \mathcal{B}(\mathbb{H})$  be the subspace of trace-class operators. For a detailed study on trace-class operators see [5, 9]. Consider an orthonormal basis  $\{w_n\}_{n \geq 1}$  for  $\mathbb{H}$ , and let

$$\mathbb{H}^1 = \left\{ f \in \mathbb{H} : \|f\| := \sum_{n=1}^{\infty} |\langle f, w_n \rangle| < \infty \right\}.$$

For a sequence  $c = (c_{mn})_{m,n=1}^{\infty} \in \ell^1$  we consider the operator  $T_c : \mathbb{H} \rightarrow \mathbb{H}$  given by

$$T_c f = \sum_{m=1}^{\infty} \sum_{n=1}^{\infty} c_{mn} \langle f, w_n \rangle w_m.$$

We say that  $T_c$  is of *Type A* with respect to the orthonormal basis  $\{w_n\}_{n \geq 1}$  if, for an orthogonal set of eigenvectors  $\{g_n\}_{n \geq 1}$  of  $T_c$  such that  $T_c = \sum_{n=1}^{\infty} g_n \otimes \overline{g_n}$ , with convergence in the strong operator topology, we have that

$$\sum_{n=1}^{\infty} \|g_n\|^2 < \infty.$$

Similarly, we say that the operator  $T_c$  is of *Type B* with respect to the orthonormal basis  $\{w_n\}_{n \geq 1}$  if there is some sequence of vectors  $\{v_n\}_{n \geq 1}$  in  $\mathbb{H}$  such that  $T_c = \sum_{n=1}^{\infty} v_n \otimes \overline{v_n}$  with convergence in the strong operator topology and we have that

$$\sum_{n=1}^{\infty} \|v_n\|^2 < \infty.$$

---

*Date:* January 10, 2020.

*2010 Mathematics Subject Classification.* Primary 45P05, 47B10; Secondary 42C15.

*Key words and phrases.* LDL factorization, Lagrangian decomposition, diagonally dominant matrices, positive semidefinite matrices, rank one matrices.

It is easy to see that if  $T_c$  is of Type  $A$  then it is of Type  $B$ . However, there exist finite rank positive trace class operators which are neither of Type  $A$  nor of Type  $B$ . We refer to [7] for more details. In [1] we proved that there exist positive trace class operators  $T_c$  of Type  $B$  which are not of Type  $A$ . Furthermore, this answers negatively a problem posed by Feichtinger [6].

Our main interest is in a finite dimensional version of the above problem. Before stating it, we set the notations that will be used through this chapter.

For  $n \geq 2$  we denote the set of all complex hermitian  $n \times n$  matrices as  $S^n := S^n(\mathbb{C})$ , positive semidefinite matrices as  $S_+^n := S_+^n(\mathbb{C})$ , and positive definite matrices  $S_{++}^n := S_{++}^n(\mathbb{C})$ . It is clear that  $S_+^n$  is a closed convex cone. Note that  $S^n = S_+^n - S_+^n$  is the (real) vector space of hermitian matrices. We will also use the notation  $U(n)$  for the set of  $n \times n$  unitary matrices.

For  $A \in S^n$ , we let  $\|A\|_{1,1} = \sum_{k,\ell=1}^n |A_{k,\ell}|$ , and we let  $\|A\|_{\mathcal{I}_1} = \sum_{k=1}^n |\lambda_k|$  where  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  are the eigenvalues of  $A$ . We recall that the operator norm of  $A \in S^n$  is given by  $\|A\|_{\text{op}} = \max\{|\lambda_k| : \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n\}$  where  $\{\lambda_k\}_{k=1}^n$  is the set of eigenvalues of  $A$ . In addition, the Frobenius norm of  $A$  is given by  $\|A\|_{\text{Fr}} = \sqrt{\text{tr}AA^*} = \sqrt{\sum_{k=1}^n \sum_{\ell=1}^n |A_{k\ell}|^2}$ . One important fact that will be used implicitly throughout the paper is that all the norms defined on  $S^n$  are equivalent and thus give rise to the same topological structure on  $S^n$ .

Similarly, for a vector  $x = (x_k)_{k=1}^n \in \mathbb{C}^n$ , and  $p \in (0, \infty)$  we let  $\|x\|_p^p = \sum_{k=1}^n |x_k|^p$  define the usual  $\ell^p$  norm,  $p \geq 1$ , with the usual modification when  $p = \infty$ , and  $p = 0$ . As pointed out above all these norms are equivalent on  $\mathbb{C}^n$  and give rise to the same topology.

The goal of this chapter is to investigate optimal decompositions of a matrix  $A \in S_+^n(\mathbb{C})$  as a sum of rank one matrices. In Section 2 we introduce some measures of optimality of the kinds of decompositions we seek, and investigate the relationship between these measures. However, before doing so, we give an exact statement of the problems we shall address and review some results about the convex cone  $S_+^n(\mathbb{C})$ . In Section 3 we restrict our attention to some classes of matrices in  $S_+^n(\mathbb{C})$ , including diagonally dominant matrices. Finally, in Section 4 we report on some numerical experiments designed to find some of these optimal decompositions.

## 2. PRELIMINARIES AND MEASURES OF OPTIMALITY

In the first part of this section, we collect some foundational facts on convex subsets of  $S^n$ . The second part will be devoted to introducing some quantities that will serve as measures of optimality of the decomposition results we seek.

**2.1. Preliminaries.** We denote the convex hull of a set  $S$  by  $\text{co}S$ . For the compact set  $X = \{xx^* : x \in \mathbb{C}^n \text{ and } \|x\|_1 = 1\}$ , we let  $\Gamma = \text{co}X$  and  $\Omega = \text{co}(X \cup \{0\})$ . Observe that  $\Omega \subset S_+^n(\mathbb{C})$ . In fact, the following result holds.

**Definition 2.1.** *An extreme point is a point such that it is not a convex combination of other points.*

**Lemma 2.2.**  *$\Omega$  is closed and compact convex subset of  $S_+^n(\mathbb{C})$  with  $\text{int } \Omega \neq \emptyset$ . Furthermore, the set of extreme points of  $\Omega$  is  $X \cup \{0\}$ .*

The proof is based on one of the versions of the Minkowski-Carathéodory Theorem, which, for completeness we recall. We refer to [3, 4, 8] for more details and background.

**Theorem 2.3.** [3, Proposition 3.1][8, Lemma 4.1] (*Minkowski-Carathéodory Theorem*) *Let  $A$  be a compact convex subset of a normed vector space  $X$  of finite dimension  $n$ . Then any point in  $A$  is a convex combination of at most  $n + 1$  extreme points. Furthermore, we can fix one of these extreme points resulting in expressing any point in  $A$  as a convex combination of at most  $n$  extreme points in addition to the one we fixed.*

*Proof of Lemma 2.2.*  $\Omega$  can be written as:

$$\begin{aligned} \Omega &= \left\{ \sum_{k=1}^m w_k x_k x_k^* : m \geq 1, \text{ an integer}, w_1, \dots, w_m \geq 0, \sum_{k=1}^m w_k \leq 1, \|x_k\|_1 = 1, 1 \leq k \leq m \right\} \\ &= \bigcup_{m \geq 1} \left\{ \sum_{k=1}^m w_k x_k x_k^* : w_1, \dots, w_m \geq 0, \sum_{k=1}^m w_k \leq 1, \|x_k\|_1 = 1, 1 \leq k \leq m \right\} \\ &= \bigcup_{m \geq 1} \Omega_m, \end{aligned}$$

where  $\Omega_m = \left\{ \sum_{k=1}^m w_k x_k x_k^* : w_1, \dots, w_m \geq 0, \sum_{k=1}^m w_k \leq 1, \|x_k\|_1 = 1, 1 \leq k \leq m \right\}$ . Notice that  $\Omega_1 \subset \Omega_2 \subset \dots \subset \Omega_m \subset \dots \subset \Omega$ . By Minkowski-Carathéodory Theorem if  $T \in \Omega$ , then

$T \in \Omega_{\dim S^n(\mathbb{C})+1}$ . Therefore

$$\begin{aligned}\Omega &= \bigcup_{m \geq 1} \Omega_m = \Omega_1 \cup \dots \cup \Omega_{n^2+1} = \Omega_{n^2+1} \\ &= \left\{ \sum_{k=1}^{n^2+1} t_k x_k x_k^* : \sum_{k=1}^{n^2+1} t_k = 1, t_k \geq 0, \|x_k\|_1 = 1, \forall k, 1 \leq k \leq n^2+1 \right\}\end{aligned}$$

We recall that the dimension of  $S^n(\mathbb{C})$  as a real vector space over is  $n^2$ . As such, and since  $X$  is compact, we conclude that  $\Omega$  as a convex hull of a compact set is compact.

To show that  $\text{int } \Omega \neq \emptyset$ , take  $\frac{1}{2n^2}I \in \Omega$ . We prove that for  $0 < r < \frac{1}{2n^2}$  we have the ball

$$B_r \left( \frac{1}{2n^2}I \right) = \left\{ \frac{1}{2n^2}I + T : T = T^*; \|T\|_{op} < r \right\} \subset \Omega.$$

Let  $T = \sum_{k=1}^n \lambda_k v_k v_k^*$ ,  $\|v_k\|_2 = 1$ , and  $|\lambda_k| \leq \|T\|_{op} < r$ . Now

$$\begin{aligned}\frac{1}{2n^2}I + T &= \frac{1}{2n^2} \sum_{k=1}^n v_k v_k^* + \sum_{k=1}^n \lambda_k v_k v_k^* \\ &= \sum_{k=1}^n \left( \frac{1}{2n^2} + \lambda_k \right) \|v_k\|_1^2 \cdot \left( \frac{v_k}{\|v_k\|_1} \right) \cdot \left( \frac{v_k}{\|v_k\|_1} \right)^*.\end{aligned}$$

Also

$$\|v_k\|_1 = \sum_{j=1}^n |v_{k,j}| \leq \left( \sum_{j=1}^n |v_{k,j}|^2 \right)^{\frac{1}{2}} \cdot \left( \sum_{j=1}^n 1 \right)^{\frac{1}{2}} = \sqrt{n} \|v_k\|_2 = \sqrt{n}.$$

Hence

$$\left\| \frac{1}{2n^2}I + T \right\|_{1,1} \leq \sum_{k=1}^n \left( \frac{1}{2n^2} + \lambda_k \right) \|v_k\|_1^2 \leq n \left( \frac{1}{2n^2} + r \right) n = \frac{1}{2} + rn^2 < 1$$

In addition, because  $r < \frac{1}{2n^2}$  we conclude that

$$\left\langle \left( \frac{1}{2n^2}I + T \right) x, x \right\rangle \geq \|x\|^2 \left( \frac{1}{2n^2} - r \right) \geq 0$$

for all  $x \in \mathbb{C}^n$ . Consequently,  $\frac{1}{2n^2}I + T \geq 0$ . We conclude that  $B_r \left( \frac{1}{2n^2}I \right) \subset \Omega$  where we use the norm  $\|A\|_{1,1}$  for convenience.  $\square$

By a similar argument,  $\Gamma$  is also compact convex subset of  $S_+^n(\mathbb{C})$ .

**2.2. Measures of optimality.** We next introduce and study the properties of some quantities defined on  $S^n$  and which will serve as measures of optimality of the rank one decompositions of matrices in  $S_+^n$ .

**Definition 2.4.** For  $A \in S_+^n$  let

$$(2.1) \quad \gamma_+(A) := \inf_{A = \sum_{n \geq 1} g_n g_n^*} \sum_{n \geq 1} \|g_n\|_1^2.$$

If  $A \in S^n$  we let

$$(2.2) \quad \gamma(A) := \inf_{A = \sum_{n \geq 1} g_n h_n^*} \sum_{n \geq 1} \|g_n\|_1 \|h_n\|_1,$$

and

$$(2.3) \quad \gamma_0(A) := \inf_{\substack{A=B-C \\ B, C \in S_+^n}} (\gamma_+(B) + \gamma_+(C)) = \inf_{A = \sum_{n \geq 1} g_n g_n^* - \sum_{k \geq 1} h_k h_k^*} \left( \sum_{n \geq 1} \|g_n\|_1^2 + \sum_{k \geq 1} \|h_k\|_1^2 \right).$$

We collect some of the properties of these functionals.

**Proposition 2.5.** *The functionals given in Definition 2.4 are sub-additive. In particular, the following statements hold.*

- (a) Given  $A, B \in S_+^n$  we have  $\gamma_+(A + B) \leq \gamma_+(A) + \gamma_+(B)$
- (b) Given  $A, B \in S^n$  we have  $\gamma(A + B) \leq \gamma(A) + \gamma(B)$
- (c) Given  $A, B \in S^n$  we have  $\gamma_0(A + B) \leq \gamma_0(A) + \gamma_0(B)$

In addition, if  $a \geq 0$ , we have  $\gamma_+(aA) = a\gamma_+(A)$  when  $A \in S_+^n$ , and

$$\begin{cases} \gamma(aA) &= |a|\gamma(A) \\ \gamma_0(aA) &= |a|\gamma_0(A) \end{cases}$$

for  $A \in S^n$  and  $a \in \mathbb{R}$ .

*Proof.* Let  $\epsilon > 0$  and choose  $\{g_k\}_{k \geq 1} \subset \mathbb{C}^n$  and  $\{h_k\}_{k \geq 1} \subset \mathbb{C}^n$  such that

$$\begin{cases} \sum_{k \geq 1} \|g_k\|_1^2 &\leq \gamma_+(A) + \epsilon/2 \\ \sum_{k \geq 1} \|h_k\|_1^2 &\leq \gamma_+(B) + \epsilon/2 \end{cases}$$

with  $A = \sum_{k \geq 1} g_k g_k^*$  and  $B = \sum_{k \geq 1} h_k h_k^*$ . It follows that

$$A + B = \sum_{k \geq 1} g_k g_k^* + \sum_{k \geq 1} h_k h_k^* = \sum_{\ell \geq 1} f_\ell f_\ell^*,$$

after reindexing. Furthermore,

$$\sum_{\ell \geq 1} \|f_\ell\|_1^2 = \sum_{k \geq 1} \|g_k\|_1^2 + \sum_{k \geq 1} \|h_k\|_1^2 \leq \gamma_+(A) + \gamma_+(B) + \epsilon.$$

The rest of the statements are proved in a similar manner, so we omit the details.  $\square$

The next result gives a comparison among the quantities defined above.

**Proposition 2.6.** *For any  $A \in S^n$  the following statements hold.*

- (a)  $\gamma(A) \leq \gamma_0(A) \leq 2\gamma(A)$ .

(b)  $\|A\|_{\mathcal{I}_1} \leq \|A\|_{1,1} \leq \gamma_0(A) \leq 2\gamma(A)$ . If in addition, we assume that  $A \in S_+^n$  then we have

$$\|A\|_{\mathcal{I}_1} \leq \|A\|_{1,1} \leq \gamma_0(A) \leq \gamma_+(A).$$

*Proof.* (a) Let  $A \in S^n$  such that  $A = A^* = \sum_{k \geq 1} g_k g_k^* - \sum_{k \geq 1} h_k h_k^*$ . Then,

$$\gamma(A) \leq \sum_{k \geq 1} \|g_k\|_1^2 + \sum_{k \geq 1} \|h_k\|_1^2.$$

Consequently,  $\gamma(A) \leq \gamma_0(A)$ .

Fix  $\varepsilon > 0$  and let  $\{g_k\}_{k=1}^M, \{h_k\}_{k=1}^M$  be such that  $A = \sum_{k=1}^M g_k h_k^*$  and

$$\sum_{k=1}^M \|g_k\|_1 \|h_k\|_1 \leq \gamma(A) + \varepsilon.$$

Furthermore, rescale  $g_k$  and  $h_k$  so that  $\|g_k\|_1 = \|h_k\|_1$ .

Let  $x_k = \frac{1}{2}(g_k + h_k)$  and  $y_k = \frac{1}{2}(g_k - h_k)$ . Then

$$\sum_{k=1}^M x_k x_k^* - \sum_{k=1}^M y_k y_k^* = \frac{1}{2} \sum_{k=1}^M g_k h_k^* + \frac{1}{2} \sum_{k=1}^M h_k g_k^* = A$$

Note also  $\|x_k\|_1 \leq \|g_k\|_1 = \|h_k\|_1$  and  $\|y_k\|_1 \leq \|g_k\|_1 = \|h_k\|_1$ . Thus

$$\gamma_0(A) \leq \sum_{k=1}^M \|x_k\|_1^2 + \sum_{k=1}^M \|y_k\|_1^2 \leq 2 \sum_{k=1}^M \|g_k\|_1^2 \leq 2\gamma(A) + 2\varepsilon.$$

Since  $\varepsilon > 0$  was arbitrary, the second inequality follows.

(b) Since  $\|A\|_{\mathcal{I}_1} = \max_{U \in U(n)} \text{Real tr}(AU)$ , let  $U_0 \in U(n)$  denote the unitary that achieves the maximum and makes the trace real. Then

$$\|A\|_{\mathcal{I}_1} = \text{tr}(AU_0) = \sum_{k=1}^n \sum_{\ell=1}^n A_{k\ell}(U_0)_{\ell k} \leq \left( \sum_{k=1}^n \sum_{\ell=1}^n |A_{k\ell}| \right) \cdot \left( \max_k \max_{\ell} |(U_0)_{\ell k}| \right) \leq \sum_{k=1}^n \sum_{\ell=1}^n |A_{k\ell}| = \|A\|_{1,1}.$$

Suppose that  $A \in S_+^n$  and let  $\varepsilon > 0$ . Choose  $\{g_k\}_{k \geq 1} \subset \mathbb{C}^n$  such that  $A = \sum_{k \geq 1} g_k g_k^*$

and

$$\sum_{k \geq 1} \|g_k\|_1^2 < \gamma_+(A) + \varepsilon.$$

It follows that

$$\gamma_0(A) \leq \sum_{k \geq 1} \|g_k\|_1^2 < \gamma_+(A) + \varepsilon.$$

□

The upper bound  $2\gamma(A)$  is tight as we show in Proposition 2.8. We next show that  $\|\cdot\|_{1,1}$  and  $\gamma(\cdot)$  are identical on  $S^n$ .

**Lemma 2.7.** *For any  $A \in S^n$  we have  $\|A\|_{1,1} = \gamma(A)$ . Consequently,  $(S^n, \gamma)$  is a normed vector space.*

*Proof.* Let  $A \in S^n$  and  $\epsilon > 0$ . Choose  $\{g_j\}_{j \geq 1}, \{h_j\}_{j \geq 1} \subset \mathbb{C}^n$  such that  $A = \sum_j g_j h_j^*$  with  $\sum_j \|g_j\|_1 \cdot \|h_j\|_1 \leq \gamma(A) + \epsilon$ . It follows that

$$\|A\|_{1,1} = \sum_{i,j} |A_{i,j}| = \left\| \sum_j g_j h_j^* \right\|_{1,1} \leq \sum_j \|g_j h_j^*\|_{1,1} \leq \sum_j \|g_j\|_1 \cdot \|h_j\|_1 \leq \gamma(A) + \epsilon.$$

Thus  $\|A\|_{1,1} \leq \gamma(A)$ .

On the other hand, for  $A \in S^n$  we can write:  $A = (A_{i,j})_{i,j} = (\sum_j (A_{i,j}))_i \cdot \delta_i^T$ , then

$$\gamma(A) \leq \sum_j \|A_{i,j}\|_1 \cdot \|\delta_i\|_1 = \sum_{i,j} |A_{i,j}| = \|A\|_{1,1}.$$

Therefore  $\|A\|_{1,1} = \gamma(A)$ . □

In fact,  $\gamma_0$  defines also a norm on  $S^n$ . More precisely, we have the following result.

**Proposition 2.8.**  *$(S^n, \gamma_0)$  is normed vector space. Furthermore,  $\gamma_0$  is Lipschitz with constant 2 on  $S^n$ :*

$$(2.4) \quad \sup_{A, B \in S^n, A \neq B} \frac{|\gamma_0(A) - \gamma_0(B)|}{\|A - B\|_{1,1}} = 2.$$

*Proof.* We have already established in Proposition 2.5 that  $\gamma_0$  satisfies the triangle inequality and is homogenous. Furthermore, suppose that  $\gamma_0(A) = 0$ . It follows that  $A = 0$ .

For the last part, let  $A, B \in S^n$ . We have

$$\gamma_0(B) = \gamma_0(B - A + A) \leq \gamma_0(B - A) + \gamma_0(A)$$

$$\gamma_0(A) = \gamma_0(B - B + A) \leq \gamma_0(B) + \gamma_0(-B + A)$$

$$\text{So } |\gamma_0(B) - \gamma_0(A)| \leq \gamma_0(B - A) \leq 2\gamma(B - A) \leq 2\|B - A\|_{1,1}.$$

To show the Lipschitz constant is exactly 2 (and hence the upper bound 2 is tight in Proposition 2.6(a) ) consider the matrix

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}.$$

Note  $\|A\|_{1,1} = 2$ . For any decomposition  $A = B - C$  with  $B, C \in S_+^2$  we have

$$B = \begin{bmatrix} a & b \\ b & c \end{bmatrix}, \quad C = \begin{bmatrix} a & e \\ e & c \end{bmatrix}$$

with  $a, c \geq 0$  and  $b - e = 1$ . Then

$$\gamma_0(A) \geq \gamma_+(B) + \gamma_+(C) \geq \gamma(B) + \gamma(C) = 2a + 2|b| + 2|1 - b| + 2c \geq 4|b| + 4|1 - b| \geq 4,$$

thanks to  $ac \geq b^2$  and  $ac \geq e^2$ . On the other hand

$$A = \frac{1}{2} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} 1 \\ -1 \end{bmatrix} \begin{bmatrix} 1 & -1 \end{bmatrix}$$

which certifies  $\gamma_0(A) = 4$ . The proof is now complete.  $\square$

We have now established that  $\gamma_0, \gamma = \|\cdot\|_{1,1}$  are equivalent norms on  $S^n$ . In addition, we proved in Proposition 2.6 that  $\gamma(A) = \|A\|_{1,1} \leq \gamma_+(A)$  for  $A \in S_+^n$ . A natural question that arises is whether a converse estimate holds. More precisely, the rest of the chapter will be devoted to investigating the following questions.

**Question 2.1.** Fix  $n \geq 2$ .

- (1) Does there exist a constant  $C > 0$ , independent of  $n$  such that for all  $A \in S_+^n$ , we have

$$\gamma_+(A) \leq C \cdot \|A\|_{1,1}.$$

- (2) For a given  $A \in S_+^n$ , give an algorithm to find  $\{h_1, h_2, \dots, h_M\}$  such that  $A = \sum_{k=1}^M h_k h_k^*$  with

$$\gamma_+(A) = \sum_{k=1}^M \|h_k\|_1^2.$$

We begin by justify why the second question makes sense. In particular, we prove that  $\gamma_+(A)$  is achieved for a certain decomposition.

**Theorem 2.9.** Given  $T \in S_+^n$ ,

$$\gamma_+(T) = \inf_{T = \sum_{k \geq 1} g_k g_k^*} \sum_{k \geq 1} \|g_k\|_1^2 = \min_{T = \sum_{k=1}^{n^2+1} g_k g_k^*} \sum_{k=1}^{n^2+1} \|g_k\|_1^2$$

for some  $\{g_k\}_{k=1}^{n^2+1} \subset \mathbb{C}^n$ .

*Proof.* Let  $T \in S_+^n(\mathbb{C})$ ,

$$\gamma_+(T) = \inf_{T = \sum_{k \geq 1} g_k g_k^*} \sum_{k \geq 1} \|g_k\|_1^2.$$

Assume  $T \neq 0$ , then  $\gamma_+(T) > 0$ . Let  $\tilde{T} = \frac{T}{\gamma_+(T)}$ ,

$$\tilde{T} = \frac{1}{\gamma_+(T)} \sum_{k \geq 1} g_k g_k^* = \sum_{k \geq 1} \frac{\|g_k\|_1^2}{\gamma_+(T)} \cdot \left( \frac{g_k}{\|g_k\|_1} \right) \cdot \left( \frac{g_k}{\|g_k\|_1} \right)^* = \sum_{k \geq 1} w_k \cdot e_k e_k^*,$$



where  $w_k = \frac{\|g_k\|_1^2}{\gamma_+(T)}$  and  $e_k = \frac{g_k}{\|g_k\|_1}$ . Hence  $\sum_{k \geq 1} w_k = \frac{1}{\gamma_+(T)} \sum_{k \geq 1} \|g_k\|_1^2 = 1$  and  $\|e_k\|_1 = 1$ . Therefore  $\gamma_+(\tilde{T}) = 1$ . It follows that  $\tilde{T} \in \Gamma$ .

By Minkowski-Carathéodory Theorem 2.3

$$\tilde{T} = \sum_{k=1}^{n^2+1} w_k \cdot e_k e_k^*, w_k \geq 0, \sum_{k=1}^{n^2+1} w_k = 1.$$

Therefore

$$\gamma_+(T) = \min_{\sum_{k=1}^{n^2+1} g_k g_k^*} \sum_{k=1}^{n^2+1} \|g_k\|_1^2.$$

□

The next question one could ask is how to find an optimal decomposition for  $A \in S_+^n$  that achieves the value  $\gamma_+(A)$ . The following technical tool will be useful in addressing this question, at least for small size matrices.

**Theorem 2.10.** *Suppose that  $A \in S_+^n(\mathbb{C})$  and  $y \in \mathbb{C}^n$ . Then  $A - yy^* \in S_+^n(\mathbb{C})$  if and only if there exists  $x \in \mathbb{C}^n$  such that  $y = Ax$  and  $\langle Ax, x \rangle \leq 1$ . When equality holds, then  $A - yy^*$  will have rank one less than that of  $A$ .*

*Proof.* The case  $y = 0$  is trivial, so we can assume without loss of generality that  $y \neq 0$ .

Suppose there exists a vector  $y$  such that  $y = Ax$  and  $\langle Ax, x \rangle \leq 1$ . For any vector  $z$  and observe that  $|\langle Ax, z \rangle|^2 \leq \langle Ax, x \rangle \langle Az, z \rangle$ . Consequently,

$$\langle (A - yy^*)z, z \rangle = \langle Az, z \rangle - |\langle Ax, z \rangle|^2 \geq \langle Az, z \rangle - \langle Ax, x \rangle \langle Az, z \rangle = \langle Az, z \rangle (1 - \langle Ax, x \rangle) \geq 0.$$

When  $\langle Ax, x \rangle = 1$ , we  $\langle (A - yy^*)x, x \rangle = \langle Ax, x \rangle - |\langle y, x \rangle|^2 = \langle Ax, x \rangle - |\langle Ax, x \rangle|^2 = 0$ . It follows that  $x \in \mathcal{N}(A - yy^*)$ . Combining the fact that  $x \notin \mathcal{N}(A)$ , we have  $\text{rank}(A - yy^*) < \text{rank}(A)$ .

For the converse, suppose that  $A - yy^*$  is positive semidefinite, where  $y \in \mathbb{C}^n$ . Write  $y = Ax + z$  where  $x \in \mathbb{C}^n$  and  $Az = 0$ . It follows that

$$\langle (A - yy^*)z, z \rangle = -|\langle y, z \rangle|^2 \leq 0$$

with equality only if  $0 = \langle z, y \rangle = \langle z, Ax + z \rangle = \langle z, z \rangle$  which implies  $z = 0$ . In addition,

$$\langle (A - yy^*)x, x \rangle = \langle Ax, x \rangle - \langle Ax, x \rangle^2 \geq 0$$

implies  $\langle Ax, x \rangle \leq 1$ .

□

The following result follows from Theorem 2.10

**Corollary 2.11.** *For any  $A \in S_+^n(\mathbb{C})$  we have*

$$\begin{aligned}\gamma_+(A) &= \min_{\langle Ax, x \rangle \leq 1, x \neq 0} \gamma_+(A - Axx^*A) + \|Ax\|_1^2 \\ &\leq \min_{\langle Ax, x \rangle = 1} \gamma_+(A - Axx^*A) + \|Ax\|_1^2.\end{aligned}$$

*Proof.* Let  $A \in S_+^n$  and  $0 \neq x \in \mathbb{C}^n$  such that  $\langle Ax, x \rangle \leq 1$ . Then by Theorem 2.10 and Proposition 2.5(a), we see that

$$\gamma_+(A) \leq \min_{\langle Ax, x \rangle \leq 1, x \neq 0} \gamma_+(A - Axx^*A) + \|Ax\|_1^2$$

On the other hand, let  $A = \sum_{k=1}^N u_k u_k^*$  be an optimal decomposition, that is  $\gamma_+(A) = \sum_{k=1}^N \|u_k\|_1^2$ . Since  $A - Axx^*A \in S_+^n$ , we can write  $A - Axx^*A = \sum_{k=1}^n v_k v_k^*$ . Hence,  $A = \sum_{k=1}^n v_k v_k^* + Axx^*A$  and by the optimality, we see that

$$\gamma_+(A - Axx^*A) + \|Ax\|_1^2 \leq \sum_{k=1}^n \|v_k\|_1^2 + \|Ax\|_1^2 \leq \gamma_+(A)$$

□

We recall that  $\Omega = \text{co}(X \cup \{0\})$  where  $X = \{xx^* : x \in \mathbb{C}^n, \|x\|_1 = 1\}$ . We now give a characterization of  $\Omega$  in terms of  $\gamma_+$  that is equivalent to the one proved in Lemma 2.2.

**Lemma 2.12.** *Using the notations of Lemma 2.2, the following result holds.  $\Omega = \{T \in S_+^n(\mathbb{C}) : \gamma_+(T) \leq 1\}$ .*

*Proof.* Let  $T \in \{T \in S_+^n(\mathbb{C}) : \gamma_+(T) \leq 1\}$ . Then

$$T = \sum_{k=1}^{n^2+1} g_k g_k^* = \sum_{k=1}^{n^2+1} w_k X_k X_k^*,$$

where  $w_k = \|g_k\|_1^2$  and  $X_k = \frac{g_k}{\|g_k\|_1}$ . Therefore  $\gamma_+(T) = \sum_{k=1}^{n^2+1} w_k \leq 1$ . Hence

$$T = \sum_{k=1}^{n^2+1} w_k X_k X_k^* + (1 - \gamma_+(T)) \cdot 0 \in \Omega.$$

Conversely, let  $T \in \Omega$ . Then  $T = \sum_k w_k X_k X_k^*$ ,  $w_k \geq 0$ , and  $\sum_k w_k \leq 1$ . Hence

$$\gamma_+(T) \leq \sum_k w_k \cdot \gamma_+(X_k X_k^*) = \sum_k w_k \leq 1.$$

□

In fact,  $\gamma_+$  can be identified with the following gauge-like function  $\varphi_\Omega : S_+^n(\mathbb{C}) \rightarrow \mathbb{R}$  defined as follows:

$$\varphi_\Omega(T) = \inf\{t > 0 : T \in t\Omega\}.$$

Let  $\tau_T = \{t > 0 : T \in t\Omega\}$ . Then  $\tau_T$  is nonempty, since  $\frac{T}{\gamma_+(T)} \in \Gamma \subset \Omega \Rightarrow T \in \gamma_+(T)\Omega \Rightarrow \gamma_+(T) \in \tau_T$ . Therefore  $\varphi_\Omega(T) \leq \gamma_+(T)$ . In fact, the following stronger result holds.

**Lemma 2.13.** *For each  $T \in S_+^n$  we have  $\varphi_\Omega(T) = \gamma_+(T)$*

*Proof.* We need to prove  $\gamma_+(T) \leq \varphi_\Omega(T)$ . If  $t \in \tau_T$ , then  $\frac{T}{t} \in \Omega$ ,

$$\begin{aligned} \frac{T}{t} &= \sum_{k=1}^{n^2+1} w_k x_k x_k^*, w_1, \dots, w_{n^2+1} \geq 0, \sum_{k=1}^{n^2+1} w_k \leq 1, \|x_k\|_1 = 1, \forall k. \\ T &= \sum_{k=1}^{n^2+1} t w_k x_k x_k^* = \sum_{k=1}^{n^2+1} g_k g_k^*, \end{aligned}$$

where  $g_k = \sqrt{t w_k} x_k$ . Now  $\gamma_+(T) \leq \sum_{k=1}^{n^2+1} t w_k = t \sum_{k=1}^{n^2+1} w_k \leq t \Rightarrow \gamma_+(T) \leq \varphi_\Omega(T)$ .  $\square$

*Remark.* It follows that  $\varphi_\Omega$  is also positively homogeneous and sub-additive, hence convex. However, we point out that  $\varphi_\Omega$  is not a Minkowski gauge function since  $\Omega$  does not include a neighborhood of 0.

We close this section with a discussion of some regularity properties of  $\gamma_+$ .

**Theorem 2.14.** *Fix  $\delta > 0$ . Let  $C_\delta = \{T \in S_+^n : T \geq \delta I, \text{tr}(T) \leq 1\}$ , then  $\gamma_+ : C_\delta \rightarrow \mathbb{R}$  is Lipschitz continuous on  $C_\delta$  with Lipschitz constant  $(n/\delta) + n^{3/2}$ .*

*Proof.* We show that,  $\forall T_1, T_2 \in C_\delta$ ,

$$|\gamma_+(T_1) - \gamma_+(T_2)| \leq \left(\frac{n}{\delta} + n^2\right) \|T_1 - T_2\|.$$

Define

$$\tilde{T} = T_2 + \frac{\delta}{\|T_2 - T_1\|} (T_2 - T_1).$$

Then

$$\lambda_{\min}(\tilde{T}) \geq \lambda_{\min}(T_2) - \left\| \frac{\delta}{\|T_2 - T_1\|} (T_2 - T_1) \right\| = \lambda_{\min}(T_2) - \delta \geq 0.$$

Consequently,  $\tilde{T} \in S_+^n$ .

Now

$$T_2 = \frac{\delta}{\delta + \|T_2 - T_1\|} T_1 + \frac{\|T_2 - T_1\|}{\delta + \|T_2 - T_1\|} \tilde{T}.$$

The convexity of  $\gamma_+$  yields

$$\gamma_+(T_2) \leq \frac{\delta}{\delta + \|T_2 - T_1\|} \gamma_+(T_1) + \frac{\|T_2 - T_1\|}{\delta + \|T_2 - T_1\|} \gamma_+(\tilde{T}),$$

which implies

$$(2.5) \quad \gamma_+(T_2) - \gamma_+(T_1) \leq \frac{\|T_2 - T_1\| (\gamma_+(\tilde{T}) - \gamma_+(T_1))}{\delta + \|T_2 - T_1\|}.$$

We have

$$(2.6) \quad \gamma_+(\tilde{T}) \leq n \cdot \text{tr}(\tilde{T}) = n \cdot \left[ \text{tr}(T_2) + \delta \cdot \text{tr} \left( \frac{T_2 - T_1}{\|T_2 - T_1\|} \right) \right] \leq n \cdot \text{tr}(T_2) + \delta n^{3/2}.$$

$$(2.7) \quad \gamma_+(T_1) \geq \|T_1\|_{1,1} = \sum_{i,j} |(T_1)_{i,j}| \geq \text{tr}(T_1) \geq n\delta.$$

Using equations (2.6) and (2.7), we get

$$(2.8) \quad \gamma_+(\tilde{T}) - \gamma_+(T_1) \leq n \cdot \text{tr}(T_2) + \delta n^{3/2} - n\delta \leq n \cdot \text{tr}(T_2) + \delta n^{3/2}.$$

Now

$$(2.9) \quad \begin{aligned} \gamma_+(T_2) - \gamma_+(T_1) &\leq \frac{\|T_2 - T_1\|}{\delta} (\gamma_+(\tilde{T}) - \gamma_+(T_1)) \leq \|T_2 - T_1\| \left[ \frac{n}{\delta} \cdot \text{tr}(T_2) + n^{3/2} \right] \\ &\Rightarrow \frac{\gamma_+(T_2) - \gamma_+(T_1)}{\|T_2 - T_1\|} \leq \frac{n}{\delta} \cdot \text{tr}(T_2) + n^{3/2}. \end{aligned}$$

Similarly

$$(2.10) \quad \frac{\gamma_+(T_1) - \gamma_+(T_2)}{\|T_1 - T_2\|} \leq \frac{n}{\delta} \cdot \text{tr}(T_1) + n^{3/2}.$$

Therefore

$$(2.11) \quad \frac{|\gamma_+(T_1) - \gamma_+(T_2)|}{\|T_1 - T_2\|} \leq \frac{n}{\delta} \cdot \max(\text{tr}(T_1), \text{tr}(T_2)) + n^{3/2} \leq \frac{n}{\delta} + n^{3/2}.$$

□

In fact, we can prove a stronger result if we restrict to  $S_{++}^n$ .

**Corollary 2.15.**  $\gamma_+ : S_{++}^n(\mathbb{C}) \rightarrow \mathbb{R}$  is continuous. Further, let  $T \in S_{++}^n(\mathbb{C})$  and  $\delta = \frac{1}{2} \lambda_{\min}(T) > 0$ . Then for every  $S \in S_{++}^n(\mathbb{C})$  with  $\|T - S\| \leq \delta$ ,

$$\frac{|\gamma_+(T) - \gamma_+(S)|}{\|T - S\|} \leq \frac{n}{\delta} \cdot \text{tr}(T) + 2n^{3/2}.$$

*Proof.* Let  $T \in S_{++}^n(\mathbb{C})$  and  $\delta = \frac{1}{2} \lambda_{\min}(T) > 0$ . For any  $S \in S_{++}^n(\mathbb{C})$  with  $\|T - S\| \leq \delta$ , and every  $x \in \mathbb{C}^n$  we have that

$$\langle Sx, x \rangle = \langle (S - T)x, x \rangle + \langle Tx, x \rangle \geq (-\delta + \lambda_{\min}(T)) \|x\|^2 = \delta \|x\|^2.$$

Using this (2.11) becomes

$$\frac{|\gamma_+(T) - \gamma_+(S)|}{\|T - S\|} \leq \frac{n}{\delta} \cdot \max(\operatorname{tr}(T), \operatorname{tr}(S)) + n^{3/2}.$$

However,  $\operatorname{tr}(S) \leq \operatorname{tr}(T) + \sqrt{n}\delta$ . Therefore,

$$\frac{|\gamma_+(T) - \gamma_+(S)|}{\|T - S\|} \leq \frac{n}{\delta} \cdot \operatorname{tr}(T) + 2n^{3/2}.$$

□

### 3. FINDING OPTIMAL RANK ONE DECOMPOSITION FOR SOME SPECIAL CLASSES OF MATRICES

In this section we consider several classes of matrices in  $S_+^n$  for which the answer to Question 2.1 is affirmative.

**3.1. Diagonally dominant matrices.** Recall that a matrix  $A \in S_+^n(\mathbb{C})$  is said to be diagonally dominant if  $A_{ii} \geq \sum_{j=1}^n |A_{ij}|$  for each  $i = 1, 2, \dots, n$ . If the inequality is strict for each  $i$ , we say that the matrix is strictly diagonally dominant. The following result can be proved for any diagonally dominant matrix in  $S_+^n$ .

**Theorem 3.1.** *Let  $A \in S_+^n$  be a diagonally dominant matrix. Then  $\gamma(A) = \gamma_0(A) = \gamma_+(A)$ .*

*Proof.* Let  $e_i = (0, \dots, 0, 1, 0, \dots, 0)$  and  $u_{ij}(x) = (0, \dots, \sqrt{x}, \dots, \overline{\sqrt{x}}, \dots, 0)$ . Given a diagonally dominant matrix  $A$ , we consider the following decomposition of  $A$  ([2])

$$A = \sum_{i < j} u_{ij}(A_{ij})u_{ij}(A_{ij})^* + \sum_i (A_{ii} - \sum_{j \in \{1, \dots, n\} \setminus \{i\}} |A_{ij}|) e_i e_i^*.$$

It follows that

$$\begin{aligned} \gamma_+(A) &\leq \sum_{i < j} 4|A_{ij}| + \sum_i (A_{ii} - \sum_{j \in \{1, \dots, n\} \setminus \{i\}} |A_{ij}|) \\ &= \sum_{i < j} 4|A_{ij}| + \sum_i A_{ii} - \sum_i \sum_{j \in \{1, \dots, n\} \setminus \{i\}} |A_{ij}| \\ &= \sum_{i < j} 4|A_{ij}| + \sum_i A_{ii} - \sum_{i < j} 2|A_{ij}| \\ &= \|A\|_{1,1}. \end{aligned}$$

□

The case of diagonally dominant matrices is a particular case of the following more general decomposition result:

**Theorem 3.2.** *Assume  $A \in S_+^n$  admits a decomposition*

$$(3.1) \quad A = \sum_{1 \leq i < j \leq n} u_{ij} u_{ij}^* + \sum_{i=1}^n v_i v_i^*$$

where each  $u_{i,j}$  has non-zero entries at most on positions  $i$  and  $j$ , and each  $v_i$  has non-zero entries at most on position  $i$ . Then  $\gamma_+(A) = \|A\|_{1,1}$ .

*Proof.* The hypothesis implies

$$u_{ij} = [0 \quad \cdots \quad 0 \quad c_{ij;i} \quad 0 \quad \cdots \quad 0 \quad c_{ij;j} \quad 0 \quad \cdots \quad 0]^T$$

and

$$v_i = [0 \quad \cdots \quad 0 \quad d_i \quad 0 \quad \cdots \quad 0]^T$$

where  $c_{ij;i}$  is on position  $i$ ,  $c_{ij;j}$  is on position  $j$  and  $d_i$  is on position  $i$ . Without loss of generality we can assume  $d_i \in \mathbb{R}$  and  $c_{ij;i}, c_{ij;j} \in \mathbb{C}$ . We write  $A = (a_{ij})_{i,j=1}^n$  where for  $1 \leq i < j \leq n$ ,  $a_{ij} = c_{ij;i} \overline{c_{ij;j}}$ , whereas for  $1 \leq i \leq n$ ,

$$a_{ii} = d_i^2 + \sum_{j=1}^{i-1} |c_{ji;i}|^2 + \sum_{j=i+1}^n |c_{ij;i}|^2.$$

These imply

$$\sum_{1 \leq i < j \leq n} \|u_{ij}\|_1^2 + \sum_{i=1}^n \|v_i\|_1^2 = \sum_{1 \leq i < j \leq n} (|u_{ij;i}| + |u_{ij;j}|)^2 + \sum_{i=1}^n d_i^2 = \sum_{1 \leq i, j \leq n} |a_{i,j}| = \|A\|_{1,1}.$$

Now the proof is complete.  $\square$

### 3.2. The cases for matrices in $S_+^n(\mathbb{C})$ for $n \in \{2, 3\}$ .

**Proposition 3.3.** *Suppose that  $A \in S_+^2$ , then*

$$\gamma_+(A) = \|A\|_{1,1}.$$

*Proof.* If  $A = uu^*$  is a rank 1 matrix in  $S_+^2$ , the proof is straightforward. Suppose  $A \in S_+^2$  is rank 2.  $A = \begin{bmatrix} a & c \\ \bar{c} & b \end{bmatrix}$  with  $ab - |c|^2 > 0$ . Using the Lagrangian decomposition [10] we can write

$$A = \begin{bmatrix} \sqrt{a} \\ \frac{c}{\sqrt{a}} \end{bmatrix} \begin{bmatrix} \sqrt{a} & \frac{c}{\sqrt{a}} \end{bmatrix} + \begin{bmatrix} 0 \\ \sqrt{b - \frac{|c|^2}{a}} \end{bmatrix} \begin{bmatrix} 0 & \sqrt{b - \frac{|c|^2}{a}} \end{bmatrix}$$

The result then follows.  $\square$

For certain  $3 \times 3$  matrices the Lagrangian decomposition [10] is optimal. In particular, we have the following result.

**Proposition 3.4.** *Let  $A \in S_+^3$  be of rank 2 or 3. If*

$$A = \begin{bmatrix} a & b & c \\ \bar{b} & d & e \\ \bar{c} & \bar{e} & f \end{bmatrix}$$

then

$$\gamma_+(A) \leq \|A\|_{1,1} + \frac{2(|ae - \bar{b}c| + |b||c| - a|e|)}{a}.$$

*In particular, if  $|ae - \bar{b}c| + |b||c| = a|e|$  then  $\gamma_+(A) = \|A\|_{1,1}$  and the Lagrangian decomposition (which in this case is the LDL factorization) is optimal.*

*Proof.* We first assume that  $A$  has rank 3. In this case,  $A$  must be positive definite and  $adf \neq 0$ . Indeed, if one of the diagonal term, say  $f = 0$ , then using the fact that  $A \in S_+^3$  would implies that  $df - |e|^2 = -|e|^2 > 0$  which is impossible.

Let

$$u_1 = \frac{1}{\sqrt{a}} A \delta_1 = \begin{bmatrix} \sqrt{a} \\ \frac{\bar{b}}{\sqrt{a}} \\ \frac{\bar{c}}{\sqrt{a}} \end{bmatrix},$$

where  $\{\delta_i\}_{i=1}^3$  is the standard ONB for  $\mathbb{C}^3$ . By Theorem 2.10, the matrix  $A - u_1 u_1^*$ . In fact, in this case, this is a rank 2 matrix given by

$$A - u_1 u_1^* = \begin{bmatrix} 0 & 0 & 0 \\ 0 & d - \frac{|b|^2}{a} & e - \frac{\bar{b}c}{a} \\ 0 & \bar{e} - \frac{\bar{c}b}{a} & f - \frac{|c|^2}{a} \end{bmatrix}$$

Let

$$u_2 = \frac{1}{\sqrt{d - \frac{|b|^2}{a}}} (A - u_1 u_1^*) \delta_2 = \begin{bmatrix} 0 \\ \sqrt{d - \frac{|b|^2}{a}} \\ \frac{\bar{e} - \frac{\bar{c}b}{a}}{\sqrt{d - \frac{|b|^2}{a}}} \end{bmatrix}.$$

It follows that  $A - u_1 u_1^* - u_2 u_2^* = u_3 u_3^*$  where

$$u_3 = \begin{bmatrix} 0 \\ 0 \\ \sqrt{\frac{\det A}{ad - |b|^2}} \end{bmatrix}.$$

Consequently, the Lagrange decomposition of  $A$  is  $A = u_1 u_1^* + u_2 u_2^* + u_3 u_3^*$  which implies that

$$\gamma_+(A) \leq \sum_{k=1}^3 \|u_k\|_1^2 = \|A\|_{1,1} + \frac{2(|ae - \bar{b}c| + |b||c| - a|e|)}{a}.$$

Now suppose that the rank of  $A$  is 2. In this case, it is possible for  $adf = 0$ . However, only one of the diagonal element can be 0. So assume that  $f = 0$ , then we also get that

$e = c = 0$ . In this case

$$A \begin{bmatrix} a & b & 0 \\ \bar{b} & d & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

which reduces to Proposition 3.3. Thus, we may assume without loss of generality that  $adf \neq 0$ . In this case, we can proceed as above. However, because the rank of the matrix  $A$  is now 2 we see that  $A = u_1 u_1^* + u_2 u_2^*$  and

$$\gamma_+(A) \leq \|u_1\|_1^2 + \|u_2\|_1^2 = \|A\|_{1,1} + \frac{2(|ae - \bar{b}c| + |b||c| - a|e|)}{a}.$$

□

*Remark.*

- (1) If one of the off diagonal elements  $b$ , or  $c$  is 0, then Proposition 3.4 shows that the Lagrange decomposition is optimal for  $\gamma_+(A)$ .

- (2) Suppose  $n = 4$  and let  $V = \frac{1}{\sqrt{14}} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & -1 \\ 1 & 1 \end{bmatrix}$ , and consider

$$A = VV^T = \frac{1}{14} \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & -1 & 1 \\ 1 & -1 & 2 & 0 \\ 1 & 1 & 0 & 2 \end{bmatrix}$$

Then  $A$  has rank 2, and the  $\|A\|_{1,1} = 1$ . However,  $\gamma_+(A) \neq \gamma(A)$ .

#### 4. NUMERICS

Here we inspect upper bounds of  $\gamma_+(A)/\|A\|_{1,1}$  for  $A$  an  $N \times N$  matrix with simulated data. We randomly generate symmetric positive definite matrices and compute upper bounds on  $\gamma_+(A)/\|A\|_{1,1}$  with different decompositions of  $A$ . The first step is generating Gaussian distributed realizations in a matrix size  $N$  by  $N$ . Then by multiplying by its transpose, the result is symmetric positive semi-definite, denoted  $A$ . Let  $\mathcal{A}_N$  denote a collection of 30 independent realizations of this random matrix.

We consider two factorizations of the matrix  $A$ : the LDL and the Eigen matrix decomposition. Specifically:

$$LDL : A = \sum_{k=1}^N v_k v_k^*$$



with  $v_k$  vectors that have the top  $k - 1$  entries 0, and

$$Eigen : A = \sum_{k=1}^n g_k g_k^*$$

where  $\{g_1, \dots, g_n\}$  are the eigenvectors, each scaled by the corresponding eigenvalue's square-root. For each decomposition denote:

$$J_{LDL}(A) = \sum_{k=1}^N \|v_k\|_1^2 \text{ and } J_{Eigen}(A) = \sum_{k=1}^N \|g_k\|_1^2$$

Let  $F_{LDL}$  and  $F_{Eigen}$  denote the worst upper bounds over the  $N$  realization ensemble:

$$F_{LDL}(N) = \max_{A \in \mathcal{A}_N} \frac{J_{LDL}(A)}{\|A\|_{1,1}}$$

$$F_{Eigen}(N) = \max_{A \in \mathcal{A}_N} \frac{J_{Eigen}(A)}{\|A\|_{1,1}}$$

We plot these worst upper bounds after 30 realizations for various  $N$  in figure 1.

In the same figure we plot the analytic approximations of these two curves using a square-root functions and a logarithmic function. The square-root function was scaled as  $c\sqrt{N}$  to closely fit the Eigen decomposition bound,  $F_{Eigen}(N)$ . Numerically we obtained  $c = 4/5$ .

From these plots we notice a clearly strictly increasing trend. Furthermore, the LDL factorization produces a smaller (tighter) upper bound than the Eigen decomposition. On the other hand, as we show in Theorem 2.9, any optimal decomposition may take  $N^2 + 1$  vectors. By limiting the number of vector to  $N$  one should not expect to achieve the optimal bound  $\gamma_+(A)$  with any decomposition.

#### ACKNOWLEDGMENTS

R. Balan was partially supported by the National Science Foundation grant DMS-1816608 and Laboratory for Telecommunication Sciences under grant H9823031D00560049. K. A. Okoudjou was partially supported by the U. S. Army Research Office grant W911NF1610008, the National Science Foundation grant DMS 1814253, and an MLK visiting professorship.

#### REFERENCES

1. R. Balan, K. A. Okoudjou, and A. Poria, *On a Feichtinger problem*, *Operators and Matrices* **12** (2018), no. 3, 881–891.
2. G. P. Barker, and D. H. Carlson, *Cones of diagonally dominant matrices*, *Pacific Jour. Math.*, bf 57 (1975), no. 1, 15–32
3. J. A. De Loera, X. Goaoc, F. Meunier, and N. H. Mustafa, *The discrete yet ubiquitous theorems of Carathéodory, Helly, Sperner, Tucker, and Tverberg*, *Bull. Amer. Math. Soc.* **56** (2019), no. 3, 415–511.
4. F. Clarke, “Functional analysis, calculus of variations and optimal control,” *Graduate Texts in Mathematics*, **264**, Springer, London, 2013.
5. N. Dunford and J. T. Schwartz, “Linear operators, Part II,” Wiley, New York, 1988.

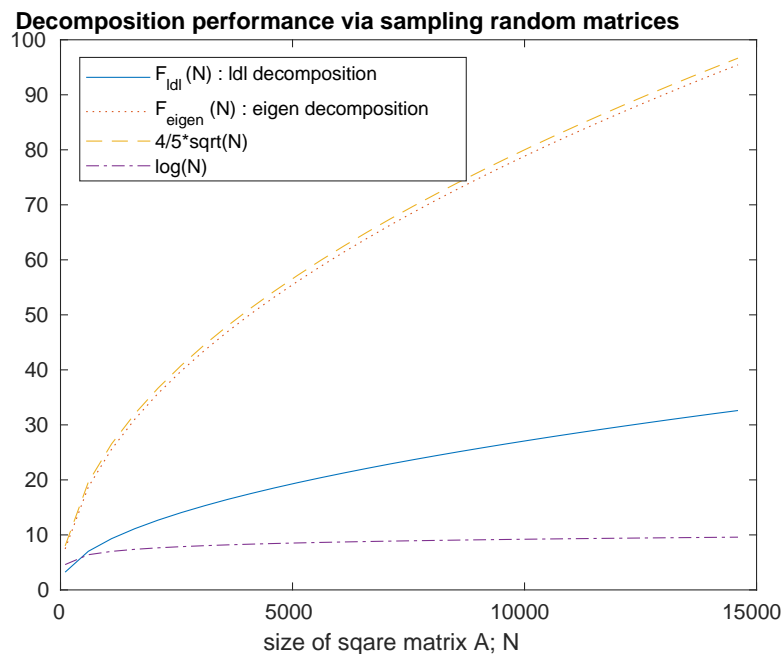


FIGURE 1

6. H. Feichtinger, P. Jorgensen, D. Larson and G. Ólafsson, *Mini-Workshop: Wavelets and Frames*, Abstracts from the mini-workshop held February 15–21, 2004, Oberwolfach Rep. **1** (2004), no. 1, 479–543.
7. C. Heil and D. Larson, *Operator theory and modulation spaces*, Contemp. Math., **451** (2008), 137–150.
8. J. Reay, *Generalizations of a theorem of Carathéodory*, Memoirs of the American Mathematical Society, vol. 54, 1965.
9. B. Simon, “Trace ideals and their applications,” Cambridge University Press, Cambridge, 1979.
10. B. Ycart, *Extreme points in convex sets of symmetric matrices*, Proc. Amer. Math. Soc., **95** (1985), no. 4, 607–612.

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF MARYLAND, COLLEGE PARK, MD 20742, USA  
 Email address: [rvbalan@umd.edu](mailto:rvbalan@umd.edu)

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF MARYLAND, COLLEGE PARK, MD 20742, USA  
 Email address: [okoudjou@umd.edu](mailto:okoudjou@umd.edu)

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF MARYLAND, COLLEGE PARK, MD 20742, USA  
 Email address: [rawson@umd.edu](mailto:rawson@umd.edu)

HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY  
 Email address: [yangwang@ust.hk](mailto:yangwang@ust.hk)

HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY  
 Email address: [zhangrui112358@yeah.net](mailto:zhangrui112358@yeah.net)