

Scalable Semidefinite Programming*

Alp Yurtsever[†], Joel A. Tropp[‡], Olivier Fercoq[§], Madeleine Udell[¶], and Volkan Cevher^{||}

Abstract. Semidefinite programming (SDP) is a powerful framework from convex optimization that has striking potential for data science applications. This paper develops a provably correct randomized algorithm for solving large, weakly constrained SDP problems by economizing on the storage and arithmetic costs. Numerical evidence shows that the method is effective for a range of applications, including relaxations of MaxCut, abstract phase retrieval, and quadratic assignment. Running on a laptop equivalent, the algorithm can handle SDP instances where the matrix variable has over 10^{14} entries.

Key words. augmented Lagrangian, conditional gradient method, convex optimization, dimension reduction, first-order method, randomized linear algebra, semidefinite programming, sketching

AMS subject classifications. 90C22, 65K05, 65F99

DOI. 10.1137/19M1305045

1. Motivation. For a spectrum of challenges in data science, methodologies based on semidefinite programming offer remarkable performance both in theory and for small problem instances. Even so, practitioners often critique this approach by asserting that it is impossible to solve semidefinite programs (SDPs) at the scale demanded by real-world applications. We would like to argue against this article of conventional wisdom.

This paper proposes a new algorithm, called SketchyCGAL, that can solve very large SDPs to moderate accuracy. The algorithm marries a primal-dual optimization technique [102] to

*Received by the editors December 9, 2019; accepted for publication (in revised form) November 16, 2020; published electronically February 8, 2021.

<https://doi.org/10.1137/19M1305045>

Funding: The first author received support from Early Postdoc.Mobility Fellowship P2ELP2 187955 from Swiss National Science Foundation and partial postdoctoral support from NSF-CAREER grant IIS-1846088. The second author was supported by ONR awards N00014-11-1-0025, N00014-17-1-2146, and N00014-18-1-2363. The fourth author was supported by DARPA award FA8750-17-2-0101. The first and fifth authors have received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program under grant agreement 725594 (time-data) and the Swiss National Science Foundation (SNSF) under grant 200021 178865/1.

[†]Laboratory for Information and Inference Systems, Department of Electrical Engineering, École Polytechnique Fédérale de Lausanne, Lausanne 1015, Switzerland, and Laboratory for Information and Decision Systems, Massachusetts Institute of Technology, Cambridge, MA 02139 USA (alp.yurtsever@epfl.ch, alpy@mit.edu, <https://lions.epfl.ch/>).

[‡]Department of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA 91125-5000 USA (jtropp@cms.caltech.edu, <http://users.cms.caltech.edu/~jtropp>).

[§]Laboratoire Traitement et Communication d’Information, Télécom Paris, Institut Polytechnique de Paris, Palaiseau 91120, France (olivier.fercoq@telecom-paris.fr, <https://perso.telecom-paristech.fr/ofercoq/>).

[¶]Department of Operations Research and Information Engineering, Cornell University, Ithaca, NY 14850 USA (udell@cornell.edu).

^{||}Laboratory for Information and Inference Systems, Department of Electrical Engineering, École Polytechnique Fédérale de Lausanne, Lausanne 1015, Switzerland (volkan.cevher@epfl.ch)

a randomized sketch for low-rank matrix approximation [93]. In each iteration, the primary expense is one low-precision randomized eigenvector calculation [57].

For every standard-form SDP that satisfies strong duality, SketchyCGAL provably converges to a near-optimal low-rank approximation of a solution. The algorithm uses limited arithmetic and minimal storage. It is most effective for weakly constrained problems whose solutions are nearly low rank. In contrast, given the same computational resources, other methods for this class of problems may fail. In particular, SketchyCGAL needs far less storage than the Burer–Monteiro factorization heuristic [24, 20] for certain problem instances [97].

In addition to the theoretical guarantees, we offer evidence that SketchyCGAL is a practical optimization algorithm. For example, on a laptop equivalent, we can solve the MaxCut SDP for a sparse graph with over 20 million vertices, where the matrix variable has over 10^{14} entries. We also tackle large phase retrieval problems arising from Fourier ptychography [48], as well as relaxations [108, 23] of the quadratic assignment problem (QAP).

1.1. Example: The maximum cut in a graph. To begin, we derive a fundamental SDP [34, 35, 44] that arises in combinatorial optimization. This example highlights why large SDPs are hard to solve, and it illustrates the potential of our approach.

1.1.1. MaxCut. Consider an undirected graph $G = (V, E)$ comprising a vertex set $V = \{1, \dots, n\}$ and a set E of m edges. The combinatorial Laplacian of the graph is the real positive-semidefinite (psd) matrix

$$(1.1) \quad L := \sum_{\{i,j\} \in E} (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^* \in \mathbb{R}^{n \times n},$$

where $\mathbf{e}_i \in \mathbb{R}^n$ denotes the i th standard basis vector and $*$ refers to the (conjugate) transpose of a matrix or vector. We can search for a maximum-weight cut in the graph by solving

$$(1.2) \quad \text{maximize } \boldsymbol{\chi}^* L \boldsymbol{\chi} \quad \text{subject to } \boldsymbol{\chi} \in \{\pm 1\}^n.$$

Unfortunately, the formulation (1.2) is NP-hard [55]. One remedy is to relax it to an SDP.

Consider the matrix $\mathbf{X} = \boldsymbol{\chi} \boldsymbol{\chi}^*$ where $\boldsymbol{\chi} \in \{\pm 1\}^n$. The matrix \mathbf{X} is psd, its diagonal entries equal one, and it has rank one. We can express the MaxCut problem (1.2) in terms of the matrix \mathbf{X} by rewriting the objective as a trace. Bringing forward the implicit constraints on \mathbf{X} and dropping the rank constraint, we arrive at the MaxCut SDP:

$$(1.3) \quad \text{maximize } \text{tr}(L\mathbf{X}) \quad \text{subject to } \text{diag}(\mathbf{X}) = \mathbf{1}, \quad \mathbf{X} \text{ is psd.}$$

As usual, diag extracts the diagonal of a matrix as a vector, and $\mathbf{1} \in \mathbb{R}^n$ is the vector of ones.

The matrix solution \mathbf{X}_* of (1.3) does not immediately yield a cut. Let $\mathbf{x}_* \mathbf{x}_*^*$ be a best rank-one approximation of \mathbf{X}_* with respect to the Frobenius norm. Then the vector $\boldsymbol{\chi}_* = \text{sgn}(\mathbf{x}_*)$ is a valid cut. In many cases, the cut $\boldsymbol{\chi}_*$ yields an excellent solution to the discrete MaxCut problem (1.2). We can also use \mathbf{X}_* to compute a cut that is provably near-optimal via a more involved randomized rounding procedure [44].

1.1.2. What's the issue? We specify an instance of the MaxCut SDP (1.3) by means of the Laplacian L of the graph, which has $O(m+n)$ nonzero entries. Our goal is to compute a best rank-one approximation of a solution to the SDP, which has $O(n)$ degrees of freedom. In other words, the total cost of representing the input and output of the problem is $O(m+n)$. Sadly, the matrix variable in (1.3) seems to require storage $O(n^2)$. For example, a graph G with one million vertices leads to an SDP (1.3) with a trillion real variables.

Storage is one of the main reasons that it has been challenging to solve large instances of the MaxCut SDP reliably. Undeterred, we raise a question:

Can we provably find a best rank-one approximation of a solution to the MaxCut SDP (1.3) with storage $O(m + n)$? Can we achieve working storage $O(n)$?

We are not aware of any correct algorithm that can solve an arbitrary instance of (1.3) with a working storage guarantee better than $\Theta(\min\{m, n^{3/2}\})$; see section 8.

In addition to the limit on storage, a good algorithm should interact with the Laplacian L only through noninvasive, low-cost operations, such as matrix-vector multiplication.

1.1.3. A storage-optimal algorithm for the MaxCut SDP. Surprisingly, it is possible to achieve all the goals announced in the last subsection.

Theorem 1.1 (MaxCut via SketchyCGAL). *For any $\varepsilon, \zeta > 0$ and any Laplacian $L \in \mathbb{R}^{n \times n}$, the SketchyCGAL algorithm computes a $(1 + \zeta)$ -optimal rank-one approximation of an ε -optimal point of (1.3); see subsection 2.2. The working storage is $O(n/\zeta)$. The algorithm performs at most $\tilde{O}(\varepsilon^{-2.5})$ matrix-vector multiplies with the Laplacian L , plus lower-order arithmetic. The algorithm is randomized; it succeeds with high probability over its random choices.*

Theorem 1.1 follows from Theorem 6.3. As usual, \tilde{O} suppresses constants and logarithmic factors. In contrast to Burer–Monteiro methods [97] and to the approximate complementarity paradigm [36], SketchyCGAL provably succeeds for every instance of MaxCut.

In our experience, SketchyCGAL works *better* than the theorem says. Figure 1.1 compares its scalability with four standard SDP solvers on a laptop equivalent (details are in subsection 7.2).

1.2. A model problem. SketchyCGAL can solve all standard-form SDPs that satisfy strong duality. To simplify parts of the presentation, we focus on a model problem that includes an extra trace constraint. Section SM4 extends SketchyCGAL to a more expressive problem template that includes standard-form SDPs with additional (conic) inequality constraints.

1.2.1. The trace-constrained SDP. We work over the field $\mathbb{F} = \mathbb{R}$ or $\mathbb{F} = \mathbb{C}$. For each $n \in \mathbb{N}$, define the set $\mathcal{S}_n := \mathcal{S}_n(\mathbb{F})$ of (conjugate) symmetric $n \times n$ matrices with entries in \mathbb{F} .

Introduce the set of $n \times n$ psd matrices with trace one:

$$(1.4) \quad \Delta_n := \{X \in \mathcal{S}_n : \text{tr } X = 1 \text{ and } X \text{ is psd}\}.$$

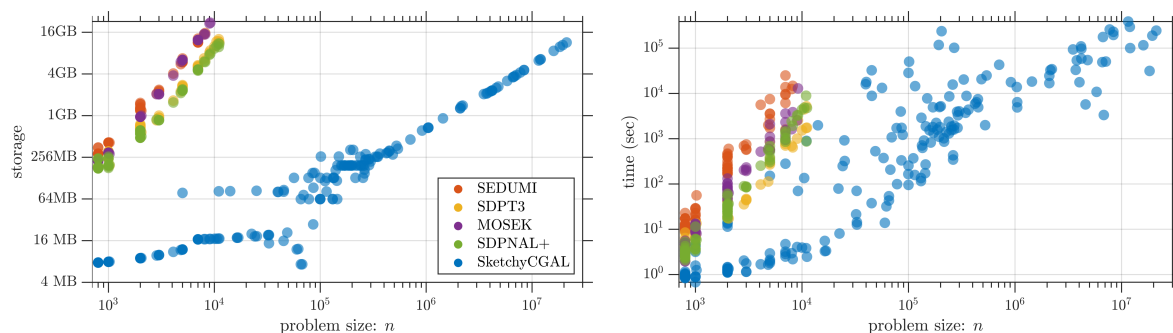


Figure 1.1. MaxCut SDP: Scalability. Storage cost (left) and runtime (right) of SketchyCGAL with sketch size $R = 10$ as compared with four standard SDP solvers. The relative error tolerance for each solver is 10^{-1} . Each marker represents one dataset. See subsection 7.2.3 for details.

Our model problem is the following trace-constrained SDP:

$$(1.5) \quad \begin{array}{ll} \text{minimize} & \text{tr}(\mathbf{C}\mathbf{X}) \\ \text{subject to} & \text{tr}(\mathbf{A}_i\mathbf{X}) = b_i \quad \text{for } i = 1, \dots, d \quad \text{and} \quad \mathbf{X} \in \alpha\Delta_n. \end{array}$$

The trace parameter $\alpha > 0$, each matrix $\mathbf{C}, \mathbf{A}_1, \dots, \mathbf{A}_d \in \mathbb{S}_n$, and $b_1, \dots, b_d \in \mathbb{R}$. We always assume that (1.5) satisfies strong duality with its standard-form dual problem.

To solve a general standard-form SDP, we replace the inclusion $\mathbf{X} \in \alpha\Delta_n$ with the constraints that \mathbf{X} is psd and $\text{tr}(\mathbf{X}) \leq \alpha$ for a large enough parameter α .

1.2.2. Applications. The model problem (1.5) has diverse applications in statistics, signal processing, quantum information theory, combinatorics, and beyond. Evidently, the MaxCut SDP (1.3) is a special case. The template (1.5) includes problems in computer vision [49], in microscopy [48], and in robotics [84]. It also supports contemporary machine learning tasks, such as certifying robustness of neural networks [82]. There is a galaxy of other examples.

1.3. Complexity of SDP formulations and solutions. This section describes some special features that commonly appear in large, real-world SDPs. Our algorithm will take advantage of these features, even as it provides guarantees for every instance of (1.5).

1.3.1. Structure of the problem data. The matrices \mathbf{C} and \mathbf{A}_i that appear in (1.5) are often highly structured, or sparse, or have low rank. As such, we can specify the SDP using a small amount of information. In our work, we exploit this property by treating the problem data for the SDP (1.5) as a collection of black boxes that support specific linear algebraic operations. The algorithm for solving the SDP only needs to access the data via these black boxes, and we can insist that these subroutines are implemented efficiently.

1.3.2. Low-rank solutions of SDPs. We will also capitalize on the fact that SDPs frequently have low-rank solutions, or the solutions are approximated well by low-rank matrices. There are several reasons why we can make this surmise.

Weakly constrained SDPs. First, many SDPs have low-rank solutions just because they are weakly constrained. That is, the number d of linear equalities in (1.5) is much smaller than the number n^2 of components in the matrix variable. This situation often occurs in signal processing and statistics problems, where d reflects the amount of measured data. SketchyCGAL is designed for weakly constrained SDPs, but it does not require this property.

A weakly constrained SDP has at least one low-rank solution because of the geometry of the set of psd matrices; see [16, Prop. II(13.4) and Prob. II.14.5] and [79].

Fact 1.2 (Barvinok–Pataki). *When $\mathbb{F} = \mathbb{R}$, the SDP (1.5) has a solution with rank $r \leq \sqrt{2(d+1)}$. When $\mathbb{F} = \mathbb{C}$, there is a solution with rank $r \leq \sqrt{d+1}$.*

For example, the MaxCut SDP (1.3) admits a solution with rank $\sqrt{2(n+1)}$.

Although a weakly constrained SDP can have solutions with high rank, a *generic* weakly constrained SDP has a unique solution, which must be low rank [5].

Fact 1.3 (Alizadeh, Haerberly, and Overton [5]). *Let $\mathbb{F} = \mathbb{R}$. Except for a set of matrices $\{\mathbf{C}, \mathbf{A}_1, \dots, \mathbf{A}_d\}$ with measure zero, the solution set of the SDP (1.5) is a unique matrix with rank $r \leq \sqrt{2(d+1)}$.*

Matrix rank minimization. Second, some SDPs are *designed* to produce a low-rank matrix that satisfies a system of linear matrix equations. This idea can be traced to the control theory literature [68, 78], and it was explored thoroughly in Fazel’s thesis [37]. Early applications include Euclidean distance matrix completion [2] and collaborative filtering [89]. Extensive empirical work indicates that these SDPs often produce low-rank solutions.

Structural properties. There are other reasons that an SDP must have a low-rank solution. For instance, consider the optimal power flow SDPs developed by Lavaei and Low [60], where the rank of the solution is controlled by the geometry of the power grid.

1.3.3. Algorithms? To summarize, many realistic SDPs have structured data, and they admit solutions that are (close to) low rank. Are there algorithms that can exploit these features? Although there are a number of methods that attempt to do so, none can provably solve every SDP with limited arithmetic and minimal storage. See [section 8](#) for related work.

Why has it been so difficult to develop provably correct algorithms for finding low-rank solutions to structured SDPs? Most approaches that try to control the rank run headlong into a computational complexity barrier: for any fixed rank parameter r , it is NP-hard to solve the model problem (1.5) if the variable \mathbf{X} is also constrained to be a rank- r matrix [37, p. 7].

To escape this sticky fact, we revise the computational goal, following [106]. The key insight is to seek a rank- r matrix that *approximates a solution* to (1.5). See [subsection 2.2](#) for a detailed explanation. This shift in perspective opens up new algorithmic prospects.

1.4. Contributions. Inspired by [106], we derive an algorithm that harnesses the favorable properties common in large SDPs. The SketchyCGAL algorithm solves the SDP (1.5) using a primal-dual optimization method [102] developed by a subset of the authors. Each iteration requires one coarse eigenvector computation [57] and leads to a rank-one update of the psd matrix variable. Instead of storing this matrix, we maintain a compressed representation by means of a matrix sketching technique [93]. After the optimization algorithm terminates, we extract from the sketch a low-rank approximation of a solution of the SDP. This idea leads to a practical, provably correct SDP solver that economizes on storage and arithmetic.

Theorem 1.4 (the model problem via SketchyCGAL). *Assume that the model problem (1.5) satisfies strong duality. For any $\varepsilon, \zeta > 0$ and any rank parameter r , the SketchyCGAL algorithm computes a $(1 + \zeta)$ -optimal rank- r approximation of an ε -optimal point of (1.5); see [subsection 2.2](#). The storage cost is $O(d + rn/\zeta)$. Most of the arithmetic consists in $\tilde{O}(\varepsilon^{-2.5})$ matrix-vector products with each matrix $\mathbf{C}, \mathbf{A}_1, \dots, \mathbf{A}_d$ from the problem data. The algorithm succeeds with high probability.*

[Theorem 6.3](#) contains full theoretical details. Note that the storage $\Theta(d + rn)$ is the minimum possible for any primal-dual algorithm that returns a rank- r solution to (1.5). [Section 7](#) contains numerical evidence that the algorithm is effective for a range of examples.

1.5. Roadmap. [Section 2](#) presents an abstract framework for studying SDPs that exposes the challenges associated with large problems. [Section 3](#) outlines a primal-dual algorithm, called CGAL, for the model problem (1.5). [Sections 4](#) and [5](#) introduce methods from randomized linear algebra that we use to control storage and arithmetic costs. [Section 6](#) develops the SketchyCGAL algorithm, its convergence theory, and its resource usage guarantees. [Section 7](#) contains a numerical study of SketchyCGAL, and [section 8](#) covers related work.

1.6. Notation. The symbol $\|\cdot\|_F$ denotes the Frobenius norm, while $\|\cdot\|_*$ is the nuclear norm (i.e., Schatten-1). The unadorned norm $\|\cdot\|$ refers to the ℓ_2 norm of a vector, the spectral norm of a matrix, or the operator norm of a linear map from $(\mathbb{S}_n, \|\cdot\|_F)$ to $(\mathbb{R}^d, \|\cdot\|)$. We write $\langle \cdot, \cdot \rangle$ for both the ℓ_2 inner product on vectors and the trace inner product on matrices.

The map $\llbracket \mathbf{M} \rrbracket_r$ returns an r -truncated singular-value decomposition of the matrix \mathbf{M} , which is a best rank- r approximation with respect to every unitarily invariant norm [69].

We use the standard computer science interpretation of the asymptotic notation O, \tilde{O}, Θ .

2. Scalable semidefinite programming. To solve the model problem (1.5) efficiently, we need to exploit structure inherent in the problem data. This section outlines an abstract approach that directs our attention to the core computational difficulties.

2.1. Abstract form of the model problem. Let us instate compact notation for the linear constraints in the model problem (1.5). Define a linear map \mathcal{A} and its adjoint \mathcal{A}^* via

$$(2.1) \quad \begin{aligned} \mathcal{A} : \mathbb{S}_n &\rightarrow \mathbb{R}^d, & \text{where } \mathcal{A}\mathbf{X} &= [\langle \mathbf{A}_1, \mathbf{X} \rangle \ \dots \ \langle \mathbf{A}_d, \mathbf{X} \rangle]; \\ \mathcal{A}^* : \mathbb{R}^d &\rightarrow \mathbb{S}_n, & \text{where } \mathcal{A}^*\mathbf{z} &= \sum_{i=1}^d z_i \mathbf{A}_i. \end{aligned}$$

We bound the linear map with the operator norm $\|\mathcal{A}\| := \|\mathcal{A}\|_{F \rightarrow \ell_2}$. Form the vector $\mathbf{b} := (b_1, \dots, b_d) \in \mathbb{R}^d$ of constraint values. In this notation, (1.5) becomes

$$(2.2) \quad \text{minimize } \langle \mathbf{C}, \mathbf{X} \rangle \quad \text{subject to } \mathcal{A}\mathbf{X} = \mathbf{b}, \quad \mathbf{X} \in \alpha \Delta_n.$$

Problem instances are parameterized by the tuple $(\mathbf{C}, \mathcal{A}, \mathbf{b}, \alpha)$.

2.2. Approximate solutions. Let \mathbf{X}_* be a solution to the model problem (2.2). For $\varepsilon \geq 0$, we say that a matrix \mathbf{X}_ε is ε -optimal for (2.2) when

$$\|\mathcal{A}\mathbf{X}_\varepsilon - \mathbf{b}\| \leq \varepsilon \quad \text{and} \quad \langle \mathbf{C}, \mathbf{X}_\varepsilon \rangle - \langle \mathbf{C}, \mathbf{X}_* \rangle \leq \varepsilon.$$

Many optimization algorithms aim to produce such ε -optimal points; cf. section 8.

As we saw in subsection 1.3.2, there are many situations where the solutions to (2.2) have low rank or they admit accurate low-rank approximations. The ε -optimal points inherit these properties for sufficiently small ε . This insight suggests a new computational goal.

For a rank parameter r , we will seek a rank- r matrix $\widehat{\mathbf{X}}$ that approximates an ε -optimal point \mathbf{X}_ε . More precisely, for a fixed suboptimality parameter $\zeta > 0$, we want

$$(2.3) \quad \|\mathbf{X}_\varepsilon - \widehat{\mathbf{X}}\|_* \leq (1 + \zeta) \cdot \|\mathbf{X}_\varepsilon - \llbracket \mathbf{X}_\varepsilon \rrbracket_r\|_*, \quad \text{where } \text{rank } \widehat{\mathbf{X}} \leq r \text{ and } \mathbf{X}_\varepsilon \text{ is } \varepsilon\text{-optimal.}$$

Given (2.3), if the ε -optimal point \mathbf{X}_ε is close to *any* rank- r matrix, then the rank- r approximate solution $\widehat{\mathbf{X}}$ is also close to the ε -optimal point \mathbf{X}_ε . This formulation is advantageous because it is easier to compute and to store the low-rank matrix $\widehat{\mathbf{X}}$.

2.3. Black-box presentation of problem data. To develop scalable algorithms for (2.2), it is productive to hide the internal complexity of the problem instance from the algorithm [106]. To do so, we treat \mathbf{C} and \mathcal{A} as black boxes that support three primitive computations:

$$(2.4) \quad \begin{array}{|c|} \hline \textcircled{1} \quad \mathbf{u} \mapsto \mathbf{C}\mathbf{u} \\ \mathbb{R}^n \rightarrow \mathbb{R}^n \\ \hline \end{array} \quad \begin{array}{|c|} \hline \textcircled{2} \quad (\mathbf{u}, \mathbf{z}) \mapsto (\mathcal{A}^*\mathbf{z})\mathbf{u} \\ \mathbb{R}^n \times \mathbb{R}^d \rightarrow \mathbb{R}^n \\ \hline \end{array} \quad \begin{array}{|c|} \hline \textcircled{3} \quad \mathbf{u} \mapsto \mathcal{A}(\mathbf{u}\mathbf{u}^*) \\ \mathbb{R}^n \rightarrow \mathbb{R}^d \\ \hline \end{array}$$

The vectors $\mathbf{u} \in \mathbb{R}^n$ and $\mathbf{z} \in \mathbb{R}^d$ are arbitrary. Although these functions may seem abstract, they are often quite natural and easy to implement well. For example, see [subsection 2.4](#).

We will formulate algorithms for (2.2) that interact with the problem data only through the operations (2.4). We tacitly assume that the primitives require minimal storage and arithmetic; otherwise, it may be impossible to develop a truly efficient algorithm.

2.4. Example: The MaxCut SDP. To provide a concrete example, let us summarize the meaning of the primitives and the desired storage costs for solving the MaxCut SDP (1.3).

- The primitive ❶ : $\mathbf{u} \mapsto -\mathbf{L}\mathbf{u}$. In the typical case that the Laplacian \mathbf{L} is sparse, this amounts to a sparse matrix-vector multiply.
- The primitive ❷ : $(\mathbf{u}, \mathbf{z}) \mapsto (\text{diag}^* \mathbf{z})\mathbf{u} = (z_1 u_1, \dots, z_n u_n)$.
- The primitive ❸ : $\mathbf{u} \mapsto \text{diag}(\mathbf{u}\mathbf{u}^*) = (|u_1|^2, \dots, |u_n|^2)$.
- The MaxCut SDP has n linear constraints, and we seek a rank-one approximation of the solution. Thus, we desire an algorithm that operates with $\Theta(n)$ working storage.

Note that we still need $\Theta(m)$ numbers to store the Laplacian \mathbf{L} of a generic graph with m edges. But we do not charge the optimization algorithm for this storage because the algorithm only interacts with \mathbf{L} through the primitives (2.4).

3. An algorithm for the model problem. We will develop a scalable method for the model problem (2.2) by enhancing an existing algorithm, called CGAL [102], developed by a subset of the authors. This method works well, but it is not as scalable as we would like because it lacks storage and arithmetic guarantees. This section summarizes the CGAL method and its convergence properties. Subsequent sections introduce additional ideas that we need to control resource usage, culminating with the SketchyCGAL algorithm in [section 6](#).

3.1. The augmented problem. For a parameter $\beta > 0$, we revise the problem (2.2) by introducing an extra term in the objective:

$$(3.1) \quad \text{minimize} \quad \langle \mathbf{C}, \mathbf{X} \rangle + \frac{\beta}{2} \|\mathcal{A}\mathbf{X} - \mathbf{b}\|^2 \quad \text{subject to} \quad \mathcal{A}\mathbf{X} = \mathbf{b}, \quad \mathbf{X} \in \alpha\Delta_n.$$

The original problem (2.2) and the augmented problem (3.1) share the same optimal value and optimal set. But the new formulation has several benefits: the augmented objective cooperates with the affine constraint to penalize infeasible points, and the dual of the augmented problem is smooth, whereas the dual of the original problem is not. See [18, Chap. 2] for background.

3.2. The augmented Lagrangian. We construct the Lagrangian L_β of the augmented problem (3.1) by introducing a dual variable $\mathbf{y} \in \mathbb{R}^d$ and promoting the affine constraint:

$$(3.2) \quad L_\beta(\mathbf{X}; \mathbf{y}) := \langle \mathbf{C}, \mathbf{X} \rangle + \langle \mathbf{y}, \mathcal{A}\mathbf{X} - \mathbf{b} \rangle + \frac{\beta}{2} \|\mathcal{A}\mathbf{X} - \mathbf{b}\|^2 \quad \text{for } \mathbf{X} \in \alpha\Delta_n \text{ and } \mathbf{y} \in \mathbb{R}^d.$$

For reference, note that the partial derivatives of the augmented Lagrangian L_β satisfy

$$(3.3) \quad \partial_{\mathbf{X}} L_\beta(\mathbf{X}; \mathbf{y}) = \mathbf{C} + \mathcal{A}^* \mathbf{y} + \beta \mathcal{A}^*(\mathcal{A}\mathbf{X} - \mathbf{b}) \quad \text{and} \quad \partial_{\mathbf{y}} L_\beta(\mathbf{X}; \mathbf{y}) = \mathcal{A}\mathbf{X} - \mathbf{b}.$$

We attempt to minimize the augmented Lagrangian L_β with respect to the primal variable \mathbf{X} , while we attempt to maximize with respect to the dual variable \mathbf{y} .

Algorithm 3.1 CGAL for the model problem (2.2)**Input:** Problem data for (2.2) implemented via the primitives (2.4); number T of iterations**Output:** Approximate solution \mathbf{X}_T to (2.2)

```

1 function CGAL( $T$ )
2   Scale problem data (subsection 7.1.1) ▷ [opt] Recommended!
3    $\beta_0 \leftarrow 1$  and  $K \leftarrow +\infty$  ▷ Default parameters
4    $\mathbf{X} \leftarrow \mathbf{0}_{n \times n}$  and  $\mathbf{y} \leftarrow \mathbf{0}_d$ 
5   for  $t \leftarrow 1, 2, 3, \dots, T$  do
6      $\beta \leftarrow \beta_0 \sqrt{t+1}$  and  $\eta \leftarrow 2/(t+1)$ 
7      $(\xi, \mathbf{v}) \leftarrow \text{ApproxMinEvec}(\mathbf{C} + \mathcal{A}^*(\mathbf{y} + \beta(\mathcal{A}\mathbf{X} - \mathbf{b})); q_t)$  ▷ Algorithm 4.1 with  $q_t = t^{1/4} \log n$ 
▷ Implement with primitives (2.4)!
8      $\mathbf{X} \leftarrow (1 - \eta)\mathbf{X} + \eta(\alpha \mathbf{v}\mathbf{v}^*)$ 
9      $\mathbf{y} \leftarrow \mathbf{y} + \gamma(\mathcal{A}\mathbf{X} - \mathbf{b})$  ▷ Step size  $\gamma$  satisfies (3.8)

```

3.3. The augmented Lagrangian strategy. The form of the augmented Lagrangian suggests an algorithm. We generate a sequence $\{(\mathbf{X}_t; \mathbf{y}_t)\}$ of primal-dual pairs by alternately minimizing over the primal variable \mathbf{X} and taking a gradient step in the dual variable \mathbf{y} .

1. **Initialize:** Choose $\mathbf{X}_1 \in \alpha \Delta_n$ and $\mathbf{y}_1 \in \mathbb{R}^d$.
2. **Primal step:** $\mathbf{X}_{t+1} \in \arg \min \{L_\beta(\mathbf{X}; \mathbf{y}_t) : \mathbf{X} \in \alpha \Delta_n\}$.
3. **Dual step:** $\mathbf{y}_{t+1} = \mathbf{y}_t + \beta(\mathcal{A}\mathbf{X} - \mathbf{b})$.

As we proceed, we can also increase the smoothing parameter β to make violations of the affine constraint in (3.1) more and more intolerable.

The augmented Lagrangian strategy is powerful, but it is hard to apply in this setting because of the cost of implementing the primal step, even approximately.

3.4. The CGAL iteration. The CGAL iteration [102] is related to the augmented Lagrangian paradigm, but the primal steps are inspired by the conditional gradient method (CGM) [39]. CGAL identifies an update direction for the primal variable by minimizing a linear proxy for the augmented Lagrangian. We take a small primal step in the update direction, and we improve the dual variable by taking a small gradient step. At each iteration, the smoothing parameter increases according to a fixed schedule.

This subsection outlines the steps in the CGAL iteration. Afterward, we give a priori guarantees on the convergence rate. Pseudocode for CGAL appears as Algorithm 3.1.

3.4.1. Initialization. Let $\beta_0 > 0$ be an initial smoothing parameter, and fix a (large) bound $K > 0$ on the maximum allowable size of the dual variable. We also assume that we have access to the norm $\|\mathcal{A}\|$ of the constraint matrix, or—failing that—a *lower* bound. Begin with an arbitrary choice $\mathbf{X}_1 \in \mathbb{S}_n$ for the primal variable; set the initial dual variable $\mathbf{y}_1 = \mathbf{0}$.

3.4.2. Primal updates via linear minimization. At iteration $t = 1, 2, 3, \dots$, we increase the smoothing parameter to $\beta_t := \beta_0 \sqrt{t+1}$, and we form the partial derivative of the augmented Lagrangian with respect to the primal variable at the current pair of iterates:

$$(3.4) \quad \mathbf{D}_t := \partial_{\mathbf{X}} L_{\beta_t}(\mathbf{X}_t; \mathbf{y}_t) = \mathbf{C} + \mathcal{A}^*(\mathbf{y}_t + \beta_t(\mathcal{A}\mathbf{X}_t - \mathbf{b})).$$

Then we compute an update $\mathbf{H}_t \in \alpha\Delta_t$ by finding the feasible point that is most correlated with the negative gradient $-\mathbf{D}_t$. This amounts to a linear minimization problem:

$$(3.5) \quad \begin{aligned} \mathbf{H}_t &\in \arg \min\{\langle \mathbf{D}_t, \mathbf{H} \rangle : \mathbf{H} \in \alpha\Delta_n\} \\ &= \text{convex hull}\{\alpha\mathbf{v}\mathbf{v}^* : \mathbf{v} \text{ is a unit-norm minimum eigenvector of } \mathbf{D}_t\}. \end{aligned}$$

In particular, we may take $\mathbf{H}_t = \alpha\mathbf{v}_t\mathbf{v}_t^*$ for a minimum eigenvector \mathbf{v}_t of the gradient \mathbf{D}_t . Next, update the matrix primal variable by taking a small step in the direction \mathbf{H}_t :

$$(3.6) \quad \mathbf{X}_{t+1} := (1 - \eta_t)\mathbf{X}_t + \eta_t\mathbf{H}_t \in \alpha\Delta_n, \quad \text{where } \eta_t := \frac{2}{t+1}.$$

The appeal of this approach is that it only requires a single eigenvector computation (3.5). Moreover, we need not compute the eigenvector accurately; see subsection 3.4.5 for details.

3.4.3. Dual updates via gradient ascent. Next, we update the dual variable by taking a small gradient step on the dual variable in the augmented Lagrangian:

$$(3.7) \quad \mathbf{y}_{t+1} = \mathbf{y}_t + \gamma_t \partial_{\mathbf{y}} L_{\beta_t}(\mathbf{X}_{t+1}; \mathbf{y}_t) = \mathbf{y}_t + \gamma_t(\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}).$$

The dual step size γ_t is the largest number that satisfies the conditions

$$(3.8) \quad \gamma_t \|\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}\|^2 \leq \frac{4\alpha^2\beta_0}{(t+1)^{3/2}} \|\mathcal{A}\|^2 \quad \text{and} \quad 0 \leq \gamma_t \leq \beta_0.$$

We omit the dual step if it makes the dual variable too large. More precisely, if (3.7) and (3.8) result in $\|\mathbf{y}_{t+1}\| > K$, then we set $\mathbf{y}_{t+1} = \mathbf{y}_t$ instead. This is the CGAL iteration.

3.4.4. Assessing solution quality. Given a primal-dual pair $(\mathbf{X}_t; \mathbf{y}_t)$, we can assess the quality of the primal solution to the model problem (2.2). First, note that we can simply compute the infeasibility: $\mathcal{A}\mathbf{X}_t - \mathbf{b}$. Second, we can bound the suboptimality of the primal objective value. To do so, define the *surrogate duality gap*

$$(3.9) \quad g_t(\mathbf{X}) := \max_{\mathbf{H} \in \alpha\Delta_n} \langle \mathbf{D}_t, \mathbf{X} - \mathbf{H} \rangle = \langle \mathbf{C}, \mathbf{X}_t \rangle + \langle \mathbf{y}_t + \beta_t(\mathcal{A}\mathbf{X}_t - \mathbf{b}), \mathcal{A}\mathbf{X}_t \rangle - \lambda_{\min}(\mathbf{D}_t).$$

The latter expression follows from (3.4) and (3.5). As usual, λ_{\min} is the minimum eigenvalue. The suboptimality of \mathbf{X}_t is bounded as

$$(3.10) \quad \langle \mathbf{C}, \mathbf{X}_t \rangle - \langle \mathbf{C}, \mathbf{X}_\star \rangle \leq g_t(\mathbf{X}_t) - \langle \mathbf{y}_t, \mathcal{A}\mathbf{X}_t - \mathbf{b} \rangle - \frac{1}{2}\beta_t \|\mathcal{A}\mathbf{X}_t - \mathbf{b}\|^2,$$

where \mathbf{X}_\star is a primal optimal point. See subsection SM1.7 for the derivation of (3.10).

For CGAL, the surrogate duality gap (3.9) is analogous to the Frank–Wolfe duality gap [50]. It offers a practical tool for detecting early convergence. Unfortunately, we have not been able to prove that the CGAL error bound (3.10) converges to zero as the algorithm converges.

3.4.5. Approximate eigenvector computations. A crucial fact is that the CGAL strategy provably works if we replace (3.5) with an approximate minimum eigenvector computation. At step t , suppose that we find an update direction

$$(3.11) \quad \mathbf{H}_t := \alpha\mathbf{v}_t\mathbf{v}_t^*, \quad \text{where } \mathbf{v}_t^* \mathbf{D}_t \mathbf{v}_t \leq \lambda_{\min}(\mathbf{D}_t) + \frac{1}{\sqrt{t+1}} \|\mathbf{D}_t\|,$$

for a unit vector $\mathbf{v}_t \in \mathbb{F}^n$. The matrix \mathbf{H}_t defined in (3.11) serves in place of a solution to (3.5) throughout the algorithm. We discuss solvers for the subproblem (3.11) in section 4.

3.5. Convergence guarantees for CGAL. The following result describes the convergence of the CGAL iteration. The analysis is adapted from [102, Thm. 3.1]; see section SM1 for details and a recapitulation of the proof.

Fact 3.1 (CGAL: Convergence). *Assume problem (2.2) satisfies strong duality. The CGAL iteration (subsection 3.4) with approximate eigenvector computations (3.11) yields a sequence $\{\mathbf{X}_t : t = 1, 2, 3, \dots\} \subset \alpha \Delta_n$ that satisfies*

$$(3.12) \quad \|\mathcal{A}\mathbf{X}_t - \mathbf{b}\| \leq \frac{\text{Const}}{\sqrt{t}} \quad \text{and} \quad |\langle \mathbf{C}, \mathbf{X}_t \rangle - \langle \mathbf{C}, \mathbf{X}_\star \rangle| \leq \frac{\text{Const}}{\sqrt{t}}.$$

The matrix \mathbf{X}_\star solves (2.2). The constant depends on the problem data $(\mathbf{C}, \mathcal{A}, \mathbf{b}, \alpha)$, the minimum Euclidean norm of a dual solution, and the algorithm parameters β_0 and K .

In view of (3.12), CGAL produces a primal iterate \mathbf{X}_T that is ε -optimal after $T = O(\varepsilon^{-2})$ iterations. The numerical work in [102, sect. 5] shows that CGAL finds an ε -optimal point after $T = O(\varepsilon^{-1})$ iterations for realistic problem instances. See also subsection SM5.3.

The analysis behind Fact 3.1 indicates that CGAL converges more quickly if we precondition the problem data by rescaling; see subsection 7.1.1. This step is critical in practice.

4. Approximate eigenvectors. Most of the computation in CGAL occurs in the approximate eigenvector step (3.11). This section describes our approach to this subproblem.

4.1. Krylov methods. The approximate minimum eigenvector computation (3.11) involves a matrix of the form $\mathbf{D}_t = \mathbf{C} + \mathcal{A}^* \mathbf{w}_t$. We can multiply the matrix \mathbf{D}_t by a vector using the primitives (2.4) 1–2. Krylov methods compute eigenvectors of \mathbf{D}_t by repeated matrix-vector multiplication with \mathbf{D}_t , so they are obvious tools for the subproblem (3.11).

Unfortunately, CGAL tends to generate matrices \mathbf{D}_t that have clustered eigenvalues. These matrices challenge standard eigenvector software, such as ARPACK [62], the engine behind the MATLAB command `eigs`. Instead, we retreat to a more classical technique.

4.2. A storage-optimal randomized Lanczos method. We use a nonstandard implementation of the randomized Lanczos method [57] to find an approximate eigenvector with minimal storage. This hinges on two ideas. First, we run the Lanczos iteration with a random initial vector, which supports convergence guarantees. Second, we do not store the Lanczos vectors but rather regenerate them to construct the approximate eigenvector. See Algorithm 4.1 for pseudocode. Kuczyński & Woźniakowski [57, Thm. 4.2(a)] have obtained error bounds.

Fact 4.1 (randomized Lanczos method). *Let $\mathbf{M} \in \mathbb{S}_n$. For $\varepsilon \in (0, 1]$ and $\delta \in (0, 0.5]$, the randomized Lanczos method, Algorithm 4.1, computes a unit vector $\mathbf{u} \in \mathbb{F}^n$ that satisfies*

$$\mathbf{u}^* \mathbf{M} \mathbf{u} \leq \lambda_{\min}(\mathbf{M}) + \frac{\varepsilon}{8} \|\mathbf{M}\| \quad \text{with probability at least } 1 - 2\delta$$

after $q \geq \frac{1}{2} + \varepsilon^{-1/2} \log(n/\delta^2)$ iterations. The arithmetic cost is at most q matrix-vector multiplies with \mathbf{M} and $O(qn)$ extra operations. The working storage is $O(n + q)$.

With constant probability, we solve the eigenvector problem (3.11) successfully in every iteration t of CGAL if we use $q_t = O(t^{1/4} \log(tn))$ iterations of Algorithm 4.1. In practice, we implement CGAL with $q_t = t^{1/4} \log n$. Although the Lanczos method has complicated numerical behavior in finite-precision arithmetic, it serves well as a subroutine within CGAL.

Algorithm 4.1 ApproxMinEvec via storage-optimal randomized Lanczos (subsection 4.2)

Input: Input matrix $M \in \mathbb{S}_n$ and maximum number q of iterations

Output: Approximate minimum eigenpair $(\xi, v) \in \mathbb{R} \times \mathbb{F}^n$ of the matrix M

```

1 function ApproxMinEvec( $M; q$ )
2    $v_1 \leftarrow \text{randn}(n, 1)$ 
3    $v_1 \leftarrow v_1 / \|v_1\|$  ▷ Store  $v_1$  for reuse
4   for  $i \leftarrow 1, 2, 3, \dots, \min\{q, n - 1\}$  do ▷ During loop, store only  $v_i$  and  $v_{i+1}$ 
5      $\omega_i \leftarrow \text{Re}(v_i^*(Mv_i))$ 
6      $v_{i+1} \leftarrow Mv_i - \omega_i v_i - \rho_{i-1} v_{i-1}$  ▷ Three-term Lanczos recurrence;  $\rho_0 v_0 = \mathbf{0}$ 
7      $\rho_i \leftarrow \|v_{i+1}\|$ 
8     if  $\rho_i = 0$  then break ▷ Found an invariant subspace!
9      $v_{i+1} \leftarrow v_{i+1} / \rho_i$ 
10     $T \leftarrow \text{tridiag}(\rho_{1:(i-1)}, \omega_{1:i}, \rho_{1:(i-1)})$  ▷ Form tridiagonal matrix
11     $[\xi, u] \leftarrow \text{MinEvec}(T)$  ▷ Exploit tridiagonal form [77, Chap. 7], or just use eig
12     $v \leftarrow \sum_{j=1}^i u_j v_j$  ▷ Modify lines 4–9 to regenerate  $v_2, \dots, v_i$  and form sum

```

Remark 4.2 (time-storage trade-off). We can retain the vectors v_i in Algorithm 4.1 to reduce the flop count by approximately a factor of two, but the storage increases to $O(qn)$.

5. Sketching and reconstruction of a psd matrix. The CGAL iteration generates a psd matrix that solves the model problem (2.2) via a sequence (3.6) of rank-one linear updates. To control storage costs, SketchyCGAL retains only a compressed version of the psd matrix variable X_t . This section outlines the *Nyström sketch*, an established method [45, 43, 64, 93] that can track the evolving psd matrix and then report a provably accurate low-rank approximation.

For more information about the implementation and behavior of low-rank matrix approximation from streaming data, see section SM2 and our papers [93, 94, 95].

5.1. Sketching and updates. Consider a psd input matrix $X \in \mathbb{S}_n$. Let R be a parameter that modulates the storage cost of the sketch and the quality of the matrix approximation.

To construct the Nyström sketch, we draw and fix a standard normal¹ test matrix $\Omega \in \mathbb{F}^{n \times R}$. Our summary, or *sketch*, of the matrix X takes the form

$$(5.1) \quad S = X\Omega \in \mathbb{F}^{n \times R}.$$

The sketch S supports linear rank-one updates to X . Indeed, we can track the evolution

$$(5.2) \quad \begin{aligned} X &\leftarrow (1 - \eta) X + \eta v v^* && \text{for } \eta \in [0, 1] \text{ and } v \in \mathbb{F}^n \\ \text{via } S &\leftarrow (1 - \eta) S + \eta v (v^* \Omega). \end{aligned}$$

The test matrix Ω and the sketch S require storage of $2Rn$ numbers in \mathbb{F} . The arithmetic cost of the linear update (5.2) to the sketch is $\Theta(Rn)$ numerical operations.

Remark 5.1 (structured random matrices). We can reduce storage costs by a factor of two by using a structured random matrix in place of Ω . For example, see [95, sect. 3] or [91].

¹Each entry of the matrix is an independent Gaussian random variable with mean zero and variance one. In the complex setting, the real and imaginary parts of each entry are independent standard normal variables.

Algorithm 5.1 NystromSketch implementation (see section 5)**Input:** Dimension n of input matrix, size R of sketch**Output:** Rank- R approximation $\widehat{\mathbf{X}}$ of sketched matrix in factored form $\widehat{\mathbf{X}} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$, where $\mathbf{U} \in \mathbb{F}^{n \times R}$ has orthonormal columns and $\mathbf{\Lambda} \in \mathbb{R}^{R \times R}$ is nonnegative diagonal

```

1 function NystromSketch.Init( $n, R$ )
2    $\mathbf{\Omega} \leftarrow \text{randn}(n, R)$  ▷ Draw and fix random test matrix
3    $\mathbf{S} \leftarrow \text{zeros}(n, R)$  ▷ Form sketch of zero matrix

4 function NystromSketch.RankOneUpdate( $\mathbf{v}, \eta$ ) ▷ Implements (5.2)
5    $\mathbf{S} \leftarrow (1 - \eta)\mathbf{S} + \eta\mathbf{v}(\mathbf{v}^*\mathbf{\Omega})$  ▷ Update sketch of matrix

6 function NystromSketch.Reconstruct()
7    $\sigma \leftarrow \sqrt{n} \text{eps}(\text{norm}(\mathbf{S}))$  ▷ Compute a shift parameter
8    $\mathbf{S}_\sigma \leftarrow \mathbf{S} + \sigma\mathbf{\Omega}$  ▷ Implicitly form sketch of  $\mathbf{X} + \sigma\mathbf{I}$ 
9    $\mathbf{L} \leftarrow \text{chol}(\mathbf{\Omega}^*\mathbf{S}_\sigma)$ 
10   $[\mathbf{U}, \mathbf{\Sigma}, \sim] \leftarrow \text{svd}(\mathbf{S}_\sigma/\mathbf{L})$  ▷ Dense SVD
11   $\mathbf{\Lambda} \leftarrow \max\{0, \mathbf{\Sigma}^2 - \sigma\mathbf{I}\}$  ▷ Remove shift

```

5.2. The reconstruction process. Given the test matrix $\mathbf{\Omega}$ and the sketch $\mathbf{S} = \mathbf{X}\mathbf{\Omega}$, we form a rank- R approximation $\widehat{\mathbf{X}}$ of the sketched matrix \mathbf{X} . The approximation is defined by

$$(5.3) \quad \widehat{\mathbf{X}} := \mathbf{S}(\mathbf{\Omega}^*\mathbf{S})^\dagger \mathbf{S}^* = (\mathbf{X}\mathbf{\Omega})(\mathbf{\Omega}^*\mathbf{X}\mathbf{\Omega})^\dagger (\mathbf{X}\mathbf{\Omega})^*,$$

where † is the pseudoinverse. This reconstruction is called a *Nyström approximation*. We often truncate $\widehat{\mathbf{X}}$ by replacing it with its best rank- r approximation $[\widehat{\mathbf{X}}]_r$ for some $r \leq R$.

See Algorithm 5.1 for a numerically stable implementation of the Nyström approximation (5.3). The algorithm takes $\Theta(R^2n)$ numerical operations and $\Theta(Rn)$ storage.

5.3. A priori error bounds. The Nyström approximation $\widehat{\mathbf{X}}$ yields a provably good estimate for the matrix \mathbf{X} contained in the sketch [93, Thm. 4.1].

Fact 5.2 (Nyström sketch: Error bound). Fix a psd matrix $\mathbf{X} \in \mathbb{S}_n$. Let $\mathbf{S} = \mathbf{X}\mathbf{\Omega}$ where $\mathbf{\Omega} \in \mathbb{F}^{n \times R}$ is standard normal. For each $r < R - 1$, the Nyström approximation (5.3) satisfies

$$(5.4) \quad \mathbb{E}_\Omega \|\mathbf{X} - \widehat{\mathbf{X}}\|_* \leq \left(1 + \frac{r}{R - r - 1}\right) \cdot \|\mathbf{X} - [\mathbf{X}]_r\|_*,$$

where \mathbb{E}_Ω is the expectation with respect to $\mathbf{\Omega}$. If we replace $\widehat{\mathbf{X}}$ with its rank- r truncation $[\widehat{\mathbf{X}}]_r$, the error bound (5.4) remains valid. Similar results hold with high probability.

6. Scalable semidefinite programming via SketchyCGAL. We may now present an extension of CGAL that solves the model problem (2.2) with controlled storage and computation. Our new algorithm, SketchyCGAL, enhances the CGAL iteration from subsection 3.4. Instead of storing the matrix \mathbf{X}_t in the CGAL iteration,

1. we drive the iteration with the d -dimensional primal state variable $\mathbf{z}_t := \mathcal{A}\mathbf{X}_t$;
2. we maintain a Nyström sketch of the primal iterate \mathbf{X}_t using storage $\Theta(Rn)$;
3. when the iteration is halted, say, at time T , we use the sketch to construct a rank- R approximation $\widehat{\mathbf{X}}_T$ of the implicitly computed solution \mathbf{X}_T of the model problem.

As we will see, the resulting method exhibits almost the same convergence behavior as the CGAL algorithm, but it also enjoys a strong storage guarantee (when $d \ll n^2$ and $R \ll n$).

6.1. The SketchyCGAL iteration. To develop the SketchyCGAL iteration, we start with the CGAL iteration. Then we make the substitutions $\mathbf{z}_t = \mathcal{A}\mathbf{X}_t$ and $\mathbf{h}_t = \mathcal{A}\mathbf{H}_t$ to eliminate the matrix variables. Let us summarize what happens; see subsection 6.2 for additional explanation. Algorithm 6.1 contains pseudocode with implementation recommendations.

6.1.1. Initialization. First, we choose the size R of the Nyström sketch. Then we draw and fix the random test matrix $\mathbf{\Omega} \in \mathbb{F}^{n \times R}$. Select an initial smoothing parameter β_0 and a (finite) bound K on the dual variable. Define the smoothing parameters $\beta_t := \beta_0 \sqrt{t+1}$ and the step size parameters $\eta_t := 2/(t+1)$. Initialize the primal state variable $\mathbf{z}_1 := \mathbf{0} \in \mathbb{R}^d$, the sketch $\mathbf{S}_1 := \mathbf{0} \in \mathbb{F}^{n \times R}$, and the dual variable $\mathbf{y}_1 := \mathbf{0} \in \mathbb{R}^d$. These choices correspond to the simplest initial iterate $\mathbf{X}_1 := \mathbf{0}$.

6.1.2. Primal updates. At iteration $t = 1, 2, 3, \dots$, we compute a unit-norm vector \mathbf{v}_t that is an approximate minimum eigenvector of the gradient \mathbf{D}_t of the smoothed objective:

$$(6.1) \quad \xi_t := \mathbf{v}_t^* \mathbf{D}_t \mathbf{v}_t \leq \lambda_{\min}(\mathbf{D}_t) + \frac{\|\mathbf{D}_t\|}{\sqrt{t+1}}, \quad \text{where } \mathbf{D}_t := \mathbf{C} + \mathcal{A}^*(\mathbf{y}_t + \beta_t(\mathbf{z}_t - \mathbf{b})).$$

This calculation corresponds with (3.11). Form the primal update direction $\mathbf{h}_t = \mathcal{A}(\alpha \mathbf{v}_t \mathbf{v}_t^*)$, and then update the primal state variable \mathbf{z}_t and the sketch \mathbf{S}_t :

$$(6.2) \quad \mathbf{z}_{t+1} := (1 - \eta_t)\mathbf{z}_t + \eta_t \mathbf{h}_t \quad \text{and} \quad \mathbf{S}_{t+1} := (1 - \eta_t)\mathbf{S}_t + \eta_t \alpha \mathbf{v}_t (\mathbf{v}_t^* \mathbf{\Omega}).$$

We obtain the update rule for the primal state variable \mathbf{z}_t by applying the linear map \mathcal{A} to the primal update rule (3.6). The update rule for the sketch \mathbf{S}_t follows from (5.2).

6.1.3. Dual updates. The update to the dual variable takes the form

$$(6.3) \quad \mathbf{y}_{t+1} = \mathbf{y}_t + \gamma_t(\mathbf{z}_{t+1} - \mathbf{b}),$$

where we choose the largest γ_t that satisfies the conditions

$$(6.4) \quad \gamma_t \|\mathbf{z}_{t+1} - \mathbf{b}\|^2 \leq \frac{4\alpha^2 \beta_0}{(t+1)^{3/2}} \|\mathcal{A}\|^2 \quad \text{and} \quad 0 \leq \gamma_t \leq \beta_0.$$

If needed, we set $\gamma_t = 0$ to prevent $\|\mathbf{y}_{t+1}\| > K$. This is the SketchyCGAL iteration.

6.2. Connection with CGAL. There is a tight connection between the iterates of SketchyCGAL and CGAL. Let $\mathbf{X}_1 := \mathbf{0}$. Using the vectors \mathbf{v}_t computed in (6.1), define matrices

$$(6.5) \quad \mathbf{H}_t := \alpha \mathbf{v}_t \mathbf{v}_t^* \quad \text{and} \quad \mathbf{X}_{t+1} := (1 - \eta_t)\mathbf{X}_t + \eta_t \mathbf{H}_t.$$

With these definitions, the following loop invariants are in force:

$$(6.6) \quad \mathbf{h}_t = \mathcal{A}\mathbf{H}_t \quad \text{and} \quad \mathbf{z}_t = \mathcal{A}\mathbf{X}_t \quad \text{and} \quad \mathbf{S}_t = \mathbf{X}_t \mathbf{\Omega}.$$

By comparing subsections 3.4 and 6.1, we see that the trajectory $\{(\mathbf{X}_t, \mathbf{H}_t, \mathbf{y}_t) : t = 1, 2, 3, \dots\}$ could also have been generated by running the CGAL iteration. In other words, the variables in SketchyCGAL track the variables of some invocation of CGAL and inherit their behavior. We refer to the matrices \mathbf{X}_t as the *implicit* CGAL iterates.

Algorithm 6.1 SketchyCGAL for the model problem (2.2)

Input: Problem data for (2.2) implemented via the primitives (2.4), sketch size R , number T of iterations
Output: Rank- R approximate solution to (2.2) in factored form $\widehat{\mathbf{X}}_T = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$, where $\mathbf{U} \in \mathbb{F}^{n \times R}$ has orthonormal columns and $\mathbf{\Lambda} \in \mathbb{R}^{R \times R}$ is nonnegative diagonal

```

1 function SketchyCGAL( $R; T$ )
2   Scale problem data (subsection 7.1.1)                                ▷ [opt] Recommended!
3    $\beta_0 \leftarrow 1$  and  $K \leftarrow +\infty$                             ▷ Default parameters
4   NystromSketch.Init( $n, R$ )
5    $\mathbf{z} \leftarrow \mathbf{0}_d$  and  $\mathbf{y} \leftarrow \mathbf{0}_d$ 
6   for  $t \leftarrow 1, 2, 3, \dots, T$  do
7      $\beta \leftarrow \beta_0 \sqrt{t+1}$  and  $\eta \leftarrow 2/(t+1)$ 
8      $[\xi, \mathbf{v}] \leftarrow \text{ApproxMinEvec}(\mathbf{C} + \mathcal{A}^*(\mathbf{y} + \beta(\mathbf{z} - \mathbf{b})); q_t)$ 
                                                    ▷ Algorithm 4.1 with  $q_t = t^{1/4} \log n$ 
                                                    ▷ Implement with primitives (2.4) 1!2!
9      $\mathbf{z} \leftarrow (1 - \eta)\mathbf{z} + \eta \mathcal{A}(\alpha \mathbf{v}^*)$ 
                                                    ▷ Use primitive (2.4) 3
10     $\mathbf{y} \leftarrow \mathbf{y} + \gamma(\mathbf{z} - \mathbf{b})$ 
                                                    ▷  $\gamma$  is the largest solution to (6.4)
11    NystromSketch.RankOneUpdate( $\sqrt{\alpha} \mathbf{v}, \eta$ )
12     $[\mathbf{U}, \mathbf{\Lambda}] \leftarrow \text{NystromSketch.Reconstruct}()$ 
13     $\mathbf{\Lambda} \leftarrow \mathbf{\Lambda} + (\alpha - \text{tr}(\mathbf{\Lambda}))\mathbf{I}_{R/R}$ 
                                                    ▷ [opt] Enforce trace constraint in (2.2)

```

6.2.1. Approximating the CGAL iterates. We do not have access to the implicit CGAL iterates described above. Nevertheless, the sketch permits us to approximate them! After iteration t of SketchyCGAL, we can form a rank- R approximation $\widehat{\mathbf{X}}_t$ of the implicit iterate \mathbf{X}_t by invoking the formula (5.3) with $\mathbf{S} = \mathbf{S}_t$. According to Fact 5.2, for each $r < R - 1$,

$$(6.7) \quad \mathbb{E}_{\Omega} \|\mathbf{X}_t - \widehat{\mathbf{X}}_t\|_* \leq \left(1 + \frac{r}{R - r - 1}\right) \cdot \|\mathbf{X}_t - \llbracket \mathbf{X}_t \rrbracket_r\|_*.$$

In other words, the computed approximation $\widehat{\mathbf{X}}_t$ is a good proxy for the implicit iterate \mathbf{X}_t whenever the latter matrix is well-approximated by a low-rank matrix. The same bound (6.7) holds if we replace $\widehat{\mathbf{X}}_t$ by the truncated rank- r matrix $\llbracket \widehat{\mathbf{X}}_t \rrbracket_r$.

Remark 6.1 (trace correction). The model problem (2.2) requires the matrix variable to have trace α , but the computed solution $\widehat{\mathbf{X}}_t$ rarely satisfies this constraint. Algorithm 6.1 includes an optional projection step (line 13) that corrects the trace. This step never increases the error in the Nyström approximation by more than a factor of two (subsection SM2.4). Our analysis of SketchyCGAL *does not* include the projection, but it is valuable in practice.

6.2.2. Assessing solution quality. Given quantities computed by SketchyCGAL, we can assess how well the implicit iterate \mathbf{X}_t solves the model problem (2.2). The infeasibility is just $\mathbf{z}_t - \mathbf{b}$. To bound suboptimality, we track the objective $p_t := \langle \mathbf{C}, \mathbf{X}_t \rangle$ by setting $p_1 = 0$ and using the update $p_{t+1} = (1 - \eta_t)p_t + \eta_t \alpha \mathbf{v}_t^*(\mathbf{C} \mathbf{v}_t)$. The surrogate duality gap (3.9) takes the form

$$g_t(\mathbf{X}_t) = p_t + \langle \mathbf{y}_t + \beta_t(\mathbf{z}_t - \mathbf{b}), \mathbf{z}_t \rangle - \lambda_{\min}(\mathbf{D}_t).$$

In practice, we use (6.1) to approximate the minimum eigenvalue: $\lambda_{\min}(\mathbf{D}_t) \approx \xi_t$. In light of (3.10), we can bound the suboptimality of \mathbf{X}_t via the expression

$$(6.8) \quad \langle \mathbf{C}, \mathbf{X}_t \rangle - \langle \mathbf{C}, \mathbf{X}_\star \rangle \leq g_t(\mathbf{X}_t) - \langle \mathbf{y}_t, \mathbf{z}_t - \mathbf{b} \rangle - \frac{1}{2} \beta_t \|\mathbf{z}_t - \mathbf{b}\|^2,$$

where \mathbf{X}_\star is a primal optimal point. See subsection SM3.1 for more details.

6.3. Convergence of SketchyCGAL. The implicit iterates \mathbf{X}_t converge to a solution of the model problem (2.2) at the same rate as the iterates of CGAL would. On average, the rank- R iterates $\widehat{\mathbf{X}}_t$ track the implicit iterates. The discrepancy between them depends on how well the implicit iterates are approximated by low-rank matrices. Here is a simple convergence result that reflects this intuition; see subsection SM3.2 for the easy proof and further results.

Theorem 6.2 (SketchyCGAL: Convergence). *Assume problem (2.2) satisfies strong duality, and let Ψ_* be its solution set. For each $r < R - 1$, the iterates $\widehat{\mathbf{X}}_t$ computed by SketchyCGAL (subsections 6.1 and 6.2) satisfy*

$$\limsup_{t \rightarrow \infty} \mathbb{E}_{\Omega} \text{dist}_*(\widehat{\mathbf{X}}_t, \Psi_*) \leq \left(1 + \frac{r}{R - r - 1}\right) \cdot \max_{\mathbf{X} \in \Psi_*} \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_*$$

The same bound holds if we replace $\widehat{\mathbf{X}}_t$ with the truncated approximation $\llbracket \widehat{\mathbf{X}}_t \rrbracket_r$. Here, dist_* is the nuclear-norm distance between a matrix and a set of matrices.

6.4. Resource usage. The main working variables in the SketchyCGAL iteration are the primal state $\mathbf{z}_t \in \mathbb{R}^d$, the computed eigenvector $\mathbf{v}_t \in \mathbb{R}^n$, the update direction $\mathbf{h}_t \in \mathbb{R}^d$, the test matrix $\Omega \in \mathbb{R}^{n \times R}$, the sketch $\mathbf{S}_t \in \mathbb{R}^{n \times R}$, and the dual variable $\mathbf{y}_t \in \mathbb{R}^d$. The total cost for storing these variables is $\Theta(d + Rn)$.

The arithmetic bottleneck in SketchyCGAL comes from the approximate eigenvector computation (6.1); see subsection 4.2 for resource requirements. The remaining computation takes place in the variable updates. To form the primal update direction \mathbf{h}_t , we invoke the primitive (2.4)Ⓜ. To update the primal state variable \mathbf{z}_t and the sketch \mathbf{S}_t in (6.2), we need $O(d + Rn)$ arithmetic operations. No further storage is required at this stage.

Table 1 documents the cost of performing T iterations of the SketchyCGAL method. The first column summarizes the resources consumed in the outer iteration. The second column tabulates the total resources spent to solve the eigenvalue problem (6.1) via the randomized Lanczos method (Algorithm 4.1).

In light of Fact 3.1 and subsection 6.2, we can be confident that the implicit iterate \mathbf{X}_T is ε -optimal within $T = O(\varepsilon^{-2})$ iterations. The formula (6.7) gives a priori guarantees on the quality of the approximation $\widehat{\mathbf{X}}_T$.

Table 1

Resource usage for T iterations of SketchyCGAL with sketch size R . The algorithm returns a rank- R approximation to an ε -optimal solution (subsection 2.2) to the model problem (2.2). In theory, $T \sim \varepsilon^{-2}$. In practice, $T \sim \varepsilon^{-1}$. Constants and logarithms are omitted. See subsection 6.5.

	SketchyCGAL Algorithm 6.1	With Lanczos Algorithm 4.1
Storage (floats)	$d + Rn$	n
Arithmetic (flops)	$T(d + Rn)$	$T^{5/4}n$
Primitives (calls)		
(2.4)ⓂⓂ	—	$T^{5/4}$
(2.4)Ⓜ	T	—

6.5. Theoretical performance of SketchyCGAL. We can package up this discussion in a theorem that describes the theoretical performance of the SketchyCGAL method.

Theorem 6.3 (SketchyCGAL). *Assume the model problem (2.2) satisfies strong duality. Fix a rank r , and set the sketch size $R \geq \zeta^{-1}r + 1$. After $T = O(\varepsilon^{-2})$ iterations, SketchyCGAL (subsections 6.1 and 6.2) returns a rank- r approximation $\widehat{\mathbf{X}} = \llbracket \widehat{\mathbf{X}}_T \rrbracket_r$ to an ε -optimal point of (2.2) that satisfies (2.3) with constant probability.*

Suppose we use the storage-optimal Lanczos method (Algorithm 4.1) to solve (6.1). Then the total storage cost for SketchyCGAL is $O(d + Rn)$ floats. The arithmetic requirements are $O(\varepsilon^{-2}(d + Rn) + \varepsilon^{-5/2}n \log(n/\varepsilon))$ flops, $O(\varepsilon^{-5/2} \log(n/\varepsilon))$ invocations of the primitives (2.4) 1 2, and $O(\varepsilon^{-2})$ applications of the primitive (2.4) 3.

The constants depend on the problem data $(\mathbf{C}, \mathcal{A}, \mathbf{b}, \alpha)$ and algorithm parameters (β_0, K) .

In practice, SketchyCGAL performs better than Theorem 6.3 suggests: empirically, we find that $T = O(\varepsilon^{-1})$ for real-world problem instances.

7. Numerical examples. This section showcases computational experiments that establish that SketchyCGAL is a practical method for solving large SDPs. We show that the algorithm is flexible by applying it to several classes of SDPs, and we show it is reliable by solving a large number of instances of each type. We give empirical evidence that the (implicit) iterates converge to optimality much faster than Theorem 6.3 suggests. Comparisons with other general-purpose SDP solvers demonstrate that SketchyCGAL is competitive for small SDPs, while it scales to problems that standard methods cannot handle.

7.1. Setup. All experiments are performed in MATLAB_R2018a with double-precision arithmetic. Source code is included with the supplementary material, linked from the main article webpage. To simulate the processing power of a personal laptop computer, we use a single Intel Xeon CPU E5-2630 v3, clocked at 2.40 GHz, with RAM usage capped at 16 GB.

Arithmetic costs are measured in terms of actual runtime. MATLAB does not currently offer a memory profiler, so we externally monitor the total memory allocated. We approximate the storage cost by reporting the peak value minus the storage at a MATLAB idle state.

7.1.1. Problem scaling. The bounds in the convergence theorem Fact 3.1 for the implicit iterates of SketchyCGAL depend on problem scaling. Our analysis motivates us to set

$$(7.1) \quad \|\mathbf{C}\|_F = \|\mathcal{A}\| = \alpha = 1 \quad \text{and} \quad \|\mathbf{A}_1\|_F = \|\mathbf{A}_2\|_F = \cdots = \|\mathbf{A}_d\|_F.$$

In our experiments, we enforce the scalings (7.1), except where noted.

7.1.2. Implementation. Our experiments require a variant of SketchyCGAL that can handle inequality constraints; see section SM4. We implement the algorithm with the default parameters ($\beta_0 = 1$ and $K = +\infty$). The sketch uses a Gaussian test matrix. The eigenvalue subproblem is solved via randomized Lanczos (Algorithm 4.1). We include the optional trace normalization (line 13 in Algorithm 6.1) whenever appropriate. *No other tuning is done.*

7.2. The MaxCut SDP. We begin with MaxCut SDPs (1.3). Our goal is to assess the storage and arithmetic costs of SketchyCGAL for a standard testbed. We compare with provable solvers for general SDPs: SEDUMI [90], SDPT3 [92], MOSEK [70], and SDPNAL+ [100].

7.2.1. Rounding. To extract a cut from an approximate solution to (1.3), we apply a simple rounding procedure. SketchyCGAL returns a matrix $\widehat{\mathbf{X}} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$, where $\mathbf{U} \in \mathbb{R}^{n \times R}$ has orthonormal columns. The columns of the entrywise signum, $\text{sgn}(\mathbf{U})$, are the signed indicators of R cuts. We compute the weights of all R cuts and select the largest. The other solvers return a full-dimensional solution $\widehat{\mathbf{X}}$; we compute the top R eigenvectors of $\widehat{\mathbf{X}}$ using the MATLAB command `eigs` and invoke the same rounding procedure.

7.2.2. Datasets. We consider datasets from two different benchmark groups:

1. GSET: 67 binary-valued matrices generated by an autonomous random graph generator and published online [101]. The dimension n varies from 800 to 10,000.
2. DIMACS10: This benchmark consists of 150 symmetric matrices (with n varying from 39 to 50,912,018) chosen for the 10th DIMACS Implementation Challenge [12]. We consider 148 datasets with dimension $n \leq 24,000,000$. Two problems (`rgg_n_2_24_s0` and `Europe_osm`) exceeded the 16 GB storage limit.

7.2.3. Storage and arithmetic comparisons. We run each solver for each dataset and measure the storage cost and the runtime. We invoke SketchyCGAL with rank parameter $R = 10$, and we stop the algorithm when the error bound (6.8) guarantees that the implicit iterates have both relative suboptimality and infeasibility below 10^{-1} . For other solvers, we set the relative error tolerance to 10^{-1} . See section SM5 for implementation details.

Figure 1.1 displays the results of this experiment. The conventional convex solvers do not scale beyond $n = 10^4$ due to their high storage costs. In contrast, SketchyCGAL handles problems with $n \approx 2.2 \cdot 10^7$. For nine datasets, SketchyCGAL did not reach the error tolerance within five days. Further details and results appear in subsection SM5.2.

7.2.4. Empirical convergence rates. Next, we investigate the empirical convergence of SketchyCGAL and the effect of the sketch size parameter R . We consider the MaxCut SDP (1.3) for the G67 dataset ($n = 10000$), the largest instance in GSET. We run 10^6 iterations of SketchyCGAL for each $R \in \{10, 25, 100\}$.

We use a high-accuracy solution from SDPT3 to approximate an optimal point \mathbf{X}_* of (2.2). Given a prospective solution \mathbf{X} , we compute its relative suboptimality and feasibility as

$$\text{objective residual} = \frac{|\langle \mathbf{C}, \mathbf{X} \rangle - \langle \mathbf{C}, \mathbf{X}_* \rangle|}{1 + |\langle \mathbf{C}, \mathbf{X}_* \rangle|} \quad \text{and} \quad \text{infeasibility} = \frac{\|\mathbf{A}\mathbf{X} - \mathbf{b}\|}{1 + \|\mathbf{b}\|}.$$

It is standard to increment the denominator by one to handle small values gracefully. These quantities are evaluated with respect to the original (not rescaled) problem data.

Figure 7.1 illustrates the convergence of SketchyCGAL for this problem. After t iterations, the residual and infeasibility decay like $O(t^{-1})$, which is far better than the $O(t^{-1/2})$ bound in Theorem 6.3. Similar behavior is manifest for other datasets and other problems.

Figure 7.1 also displays the weight of the cut obtained after rounding, compared with the weight of the cut obtained from the SDPT3 solution. Observe that sketch size $R = 10$ is sufficient, and the SketchyCGAL solutions yield excellent cuts after a few hundred iterations.

7.2.5. Primal-dual convergence. We have observed empirically that convergence occurs for the implicit primal sequence (\mathbf{X}_t) , the dual sequence (\mathbf{y}_t) , and the posterior error bound (6.8) generated by SketchyCGAL. See subsection SM5.3 for numerical evidence.

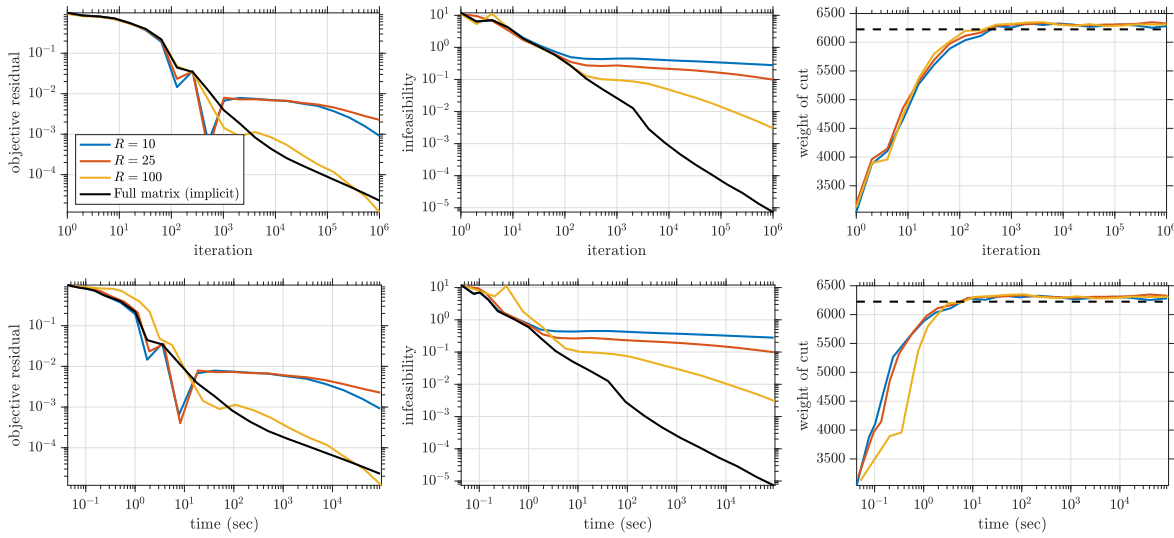


Figure 7.1. *MaxCut SDP: Convergence.* We solve the MaxCut SDP for the G67 dataset ($n = 10000$) with SketchyCGAL. The subplots show the suboptimality (left), infeasibility (center), and cut value (right) of the implicit iterates and low-rank approximations for sketch size $R \in \{10, 25, 100\}$ as a function of iteration (top) and runtime (bottom). The dashed line is the cut value of a high-accuracy SDPT3 solution. See subsection 7.2.4.

7.2.6. Hard MaxCut instances. Waldspurger and Waters [97] construct instances of the MaxCut SDP (1.3) that are challenging for algorithms based on the Burer–Monteiro [24] factorization heuristic (subsection 8.4). Each instance has a unique solution and the solution has rank 1, but Burer–Monteiro methods require factorization rank $R = \Theta(\sqrt{n})$, resulting in storage cost $\Theta(n^{3/2})$. In subsection SM5.4, we give numerical evidence that SketchyCGAL can solve these instances with sketch size $R = 2$, achieving the optimal storage $\Theta(n)$. We confirm that Burer–Monteiro usually fails in the optimal-storage regime.

7.3. Abstract phase retrieval. Phase retrieval is the problem of reconstructing a complex-valued signal from intensity-only measurements. It arises in interferometry [38], speech processing [15], array imaging [30], microscopy [48], and many other applications. We will outline a standard method [30, 48] for performing phase retrieval by means of an SDP.

This section uses synthetic instances of a phase retrieval SDP to compare the scaling behavior of SketchyCGAL and CGAL. We also consider a third algorithm ThinCGAL, inspired by [104], that maintains a thin SVD of the matrix variable via rank-one updates [22].

7.3.1. Phase retrieval SDPs. Let $\chi_{\natural} \in \mathbb{C}^n$ be an unknown (discrete) signal. For known vectors $\mathbf{a}_i \in \mathbb{C}^n$, we acquire measurements of the form

$$(7.2) \quad b_i = |\langle \mathbf{a}_i, \chi_{\natural} \rangle|^2 \quad \text{for } i = 1, 2, 3, \dots, d.$$

Abstract phase retrieval is the challenging problem of recovering χ_{\natural} from \mathbf{b} .

Let us summarize a lifting approach introduced by Balan et al. [14]. The key idea is to replace the signal vector \mathbf{x}_{\natural} by the matrix $\mathbf{X}_{\natural} = \chi_{\natural} \chi_{\natural}^*$. Then rewrite (7.2) as

$$(7.3) \quad b_i = \mathbf{a}_i^* \mathbf{X}_{\natural} \mathbf{a}_i = \langle \mathbf{A}_i, \mathbf{X}_{\natural} \rangle, \quad \text{where } \mathbf{A}_i = \mathbf{a}_i \mathbf{a}_i^* \quad \text{for } i = 1, 2, 3, \dots, d.$$

Promoting the implicit constraints on \mathbf{X}_\dagger and forming the \mathbf{A}_i into a linear map \mathcal{A} , we can express the problem of finding \mathbf{X}_\dagger as a feasibility problem with a matrix variable:

$$\text{find } \mathbf{X} \in \mathbb{S}_n \text{ subject to } \mathcal{A}\mathbf{X} = \mathbf{b}, \quad \mathbf{X} \text{ is psd}, \quad \text{rank}(\mathbf{X}) = 1.$$

To reach a tractable convex formulation, we pass to a trace minimization SDP [37]:

$$(7.4) \quad \text{minimize } \text{tr } \mathbf{X} \quad \text{subject to } \mathcal{A}\mathbf{X} = \mathbf{b}, \quad \mathbf{X} \text{ is psd}, \quad \text{tr } \mathbf{X} \leq \alpha.$$

The parameter α is an upper bound on the signal energy $\|\chi_\dagger\|^2$, which can be estimated from the observed data \mathbf{b} ; see [104] for the details. We can solve the SDP (7.4) via SketchyCGAL.

7.3.2. Rounding. Suppose that we have obtained an approximate solution \mathbf{X} to (7.4). To estimate the signal χ_\dagger , we form the vector $\chi = \sqrt{\lambda}\mathbf{u}$ where (λ, \mathbf{u}) is a maximum eigenpair of \mathbf{X} . Both SketchyCGAL and ThinCGAL return eigenvalue decompositions, so this step is trivial. For CGAL, we use the MATLAB function `eigs` to perform this computation.

7.3.3. Dataset and evaluation. We consider synthetic phase retrieval instances. For each $n \in \{10^2, 10^3, \dots, 10^6\}$, we generate 20 independent datasets as follows. First, draw $\chi_\dagger \in \mathbb{C}^n$ from the complex standard normal distribution. Then acquire $d = 12n$ phaseless measurements (7.2) using the coded diffraction pattern model [28]; see subsection SM6.1. The induced linear maps \mathcal{A} and \mathcal{A}^* can be applied via the fast Fourier transform (FFT).

The relative error in a signal reconstruction χ is given by $\min_{\phi \in \mathbb{R}} \|e^{i\phi}\chi - \chi_\dagger\| / \|\chi_\dagger\|$. In the SDP (7.4), we set $\alpha = 3n$ to demonstrate insensitivity of the algorithm to the choice of α .

7.3.4. Storage and arithmetic comparisons. For each algorithm, we report the storage cost and runtime required to produce a signal estimate χ with (exact) relative error below 10^{-2} . We invoke SketchyCGAL with sketch size parameter $R = 5$.

Figure 7.2 displays the outcome. We witness the benefit of sketching for both storage and arithmetic. CGAL fails for all large-scale instances ($n = 10^5$ and 10^6) due to storage allocation. For the same reason, ThinCGAL fails for $n = 10^6$. SketchyCGAL successfully solves all problem instances to the target accuracy.

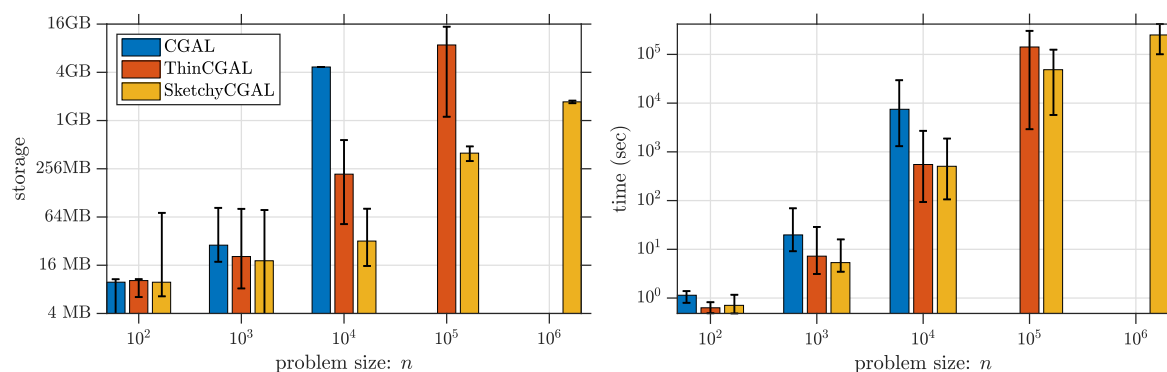


Figure 7.2. Phase retrieval SDP: Scalability. Storage cost (left) and runtime (right) to solve random instances with algorithms CGAL, ThinCGAL, and SketchyCGAL. The height of each bar is the mean; the interval marks the minimum and maximum over 20 trials. Missing bars indicate total failure. See subsection 7.3.4.

7.4. Phase retrieval in microscopy. Next, we study a more realistic phase retrieval problem that arises from a type of microscopy system [110] called Fourier ptychography (FP). Phase retrieval SDPs offer a potential approach to FP imaging [48]. This section shows that SketchyCGAL can successfully solve the difficult phase retrieval SDPs that arise from FP.

7.4.1. Fourier ptychography. FP microscopes circumvent the physical limits of a simple lens to achieve high-resolution and wide field-of-view simultaneously [110]. To do so, an FP microscope illuminates a sample from many angles and uses a simple lens to collect low-resolution intensity-only images. The measurements are low-pass filters, whose transfer functions depend on the lens and the angle of illumination [48]. From the data, we form a high-resolution image by solving a phase retrieval problem, e.g., via the SDP (7.4).

The high-resolution image of the sample is represented by a Fourier-domain vector $\chi_{\dagger} \in \mathbb{C}^n$. We acquire d intensity-only measurements of the form (7.2), where d is the total number of pixels in the low-resolution illuminations. The low-pass measurements are encoded in vectors \mathbf{a}_i . The operators \mathcal{A} and \mathcal{A}^* , built from the matrices $\mathbf{A}_i = \mathbf{a}_i \mathbf{a}_i^*$, can be applied via the FFT.

7.4.2. Dataset and evaluation. The authors of [48] provided transmission matrices \mathbf{A}_i of a working FP system. We simulate this system in the computer environment to acquire noiseless intensity-only measurements of a high-resolution target image [33]. In this setup, $\chi_{\dagger} \in \mathbb{C}^n$ corresponds to the Fourier transform of an $n = 320^2 = 102,400$ pixel grayscale image. We normalize χ_{\dagger} so that $\|\chi_{\dagger}\| = 1$. We acquire 225 low-resolution illuminations of the original image, each with $64^2 = 4,096$ pixels. The total number of measurements is $d = 921,600$.

We evaluate the error of an estimate χ as in subsection 7.3.3. Although we know that the signal energy equals 1, it is more realistic to approximate the signal energy by setting $\alpha = 1.5$ in the SDP (7.4).

7.4.3. FP imaging. We solve the phase retrieval SDP (7.4) by performing 10,000 iterations of SketchyCGAL with rank parameter $R = 5$. The top eigenvector of the output gives an approximation $\chi \in \mathbb{C}^n$ of the signal. The inverse Fourier transform of χ is the desired image.

Figure 7.3 displays the image reconstruction after $t \in \{10, 10^2, 10^3, 10^4\}$ iterations. We obtain a good-quality result in 5 hours after 1,000 iterations; a sharper image emerges in 59 hours after 10,000 iterations are complete. We believe the computational time can be reduced substantially with a parallel or GPU implementation. Nevertheless, it is gratifying that we have solved a difficult SDP whose matrix variable has over 10^{10} entries. See subsection SM6.2 for a larger FP instance with $n = 640^2$ pixels.

7.5. The quadratic assignment problem. The QAP is a very difficult combinatorial optimization problem that includes the traveling salesman, max-clique, and bandwidth problems and many others as special cases [66]. SDP relaxations offer a powerful approach for obtaining good solutions to large QAP problems [108]. In this section, we demonstrate that SketchyCGAL can solve these challenging SDPs.

7.5.1. QAP. We begin with the simplest form of the QAP. Fix symmetric $\mathbf{n} \times \mathbf{n}$ matrices $\mathbf{A}, \mathbf{B} \in \mathbb{S}_{\mathbf{n}}$ where \mathbf{n} is a natural number. We wish to “align” the matrices by solving

$$(7.5) \quad \text{minimize} \quad \text{tr}(\mathbf{A}\mathbf{\Pi}\mathbf{B}\mathbf{\Pi}^*) \quad \text{subject to} \quad \mathbf{\Pi} \text{ is an } \mathbf{n} \times \mathbf{n} \text{ permutation matrix.}$$

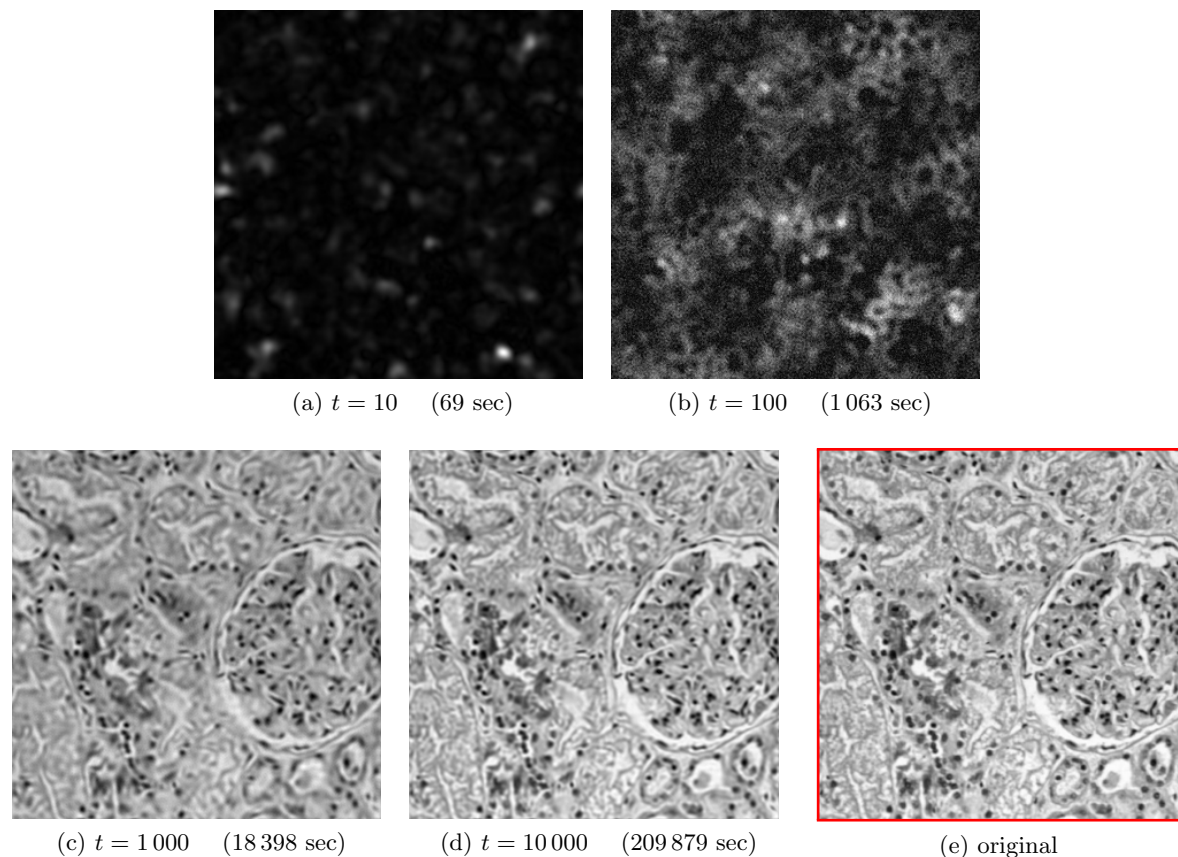


Figure 7.3. Phase retrieval SDP: Imaging. Reconstruction of an $n = 320^2$ pixel image from FP data. We solve an $n \times n$ phase retrieval SDP via SketchyCGAL with rank parameter $R = 5$ and show the images obtained at iterations $t = 10^3, 10^4$. The last subfigure is the original. See subsection 7.4.3.

Recall that a permutation matrix $\mathbf{\Pi}$ has precisely one nonzero entry in each row and column and that nonzero entry equals one. A brute force search over the $n!$ permutation matrices of size n quickly becomes intractable as n grows. Unsurprisingly, the QAP problem (7.5) is NP-hard [86]. Instances with $n > 30$ usually cannot be solved in reasonable time.

7.5.2. Relaxations. There is an extensive literature on SDP relaxations for QAP, beginning with the work [108] of Zhao et al. We consider a weaker relaxation inspired by [49, 23]:

$$\begin{aligned}
 (7.6) \quad & \text{minimize} && \text{tr}[(\mathbf{B} \otimes \mathbf{A})\mathbf{Y}] \\
 & \text{subject to} && \text{tr}_1(\mathbf{Y}) = \mathbf{I}, \quad \text{tr}_2(\mathbf{Y}) = \mathbf{I}, \quad \mathcal{G}(\mathbf{Y}) \geq 0, \\
 & && \text{vec}(\mathbf{P}) = \text{diag}(\mathbf{Y}), \quad \mathbf{P}\mathbf{1} = \mathbf{1}, \quad \mathbf{1}^*\mathbf{P} = \mathbf{1}^*, \quad \mathbf{P} \geq \mathbf{0}, \\
 & && \begin{bmatrix} 1 & \text{vec}(\mathbf{P})^* \\ \text{vec}(\mathbf{P}) & \mathbf{Y} \end{bmatrix} \succcurlyeq \mathbf{0}, \quad \text{tr } \mathbf{Y} = n.
 \end{aligned}$$

We have written \otimes for the Kronecker product.

The constraint $\mathcal{G}(\mathbf{Y}) \geq 0$ enforces nonnegativity of a subset of the entries in \mathbf{Y} . In the Zhao et al. relaxation, \mathcal{G} is the identity map, so it yields $O(\mathbf{n}^4)$ constraints. We reduce the complexity by choosing \mathcal{G} more carefully. In our formulation, \mathcal{G} extracts precisely the nonzero entries of the matrix $\mathbf{B} \otimes \mathbf{1}\mathbf{1}^*$. This is beneficial because \mathbf{B} is sparse in many applications. For example, in traveling salesman and bandwidth problems, \mathbf{B} has $O(\mathbf{n})$ nonzero entries, so the map \mathcal{G} produces only $O(\mathbf{n}^3)$ constraints.

The main variable \mathbf{Y} in (7.6) has dimension $\mathbf{n}^2 \times \mathbf{n}^2$. As a consequence, the problem has $O(\mathbf{n}^4)$ degrees of freedom, together with $O(\mathbf{n}^3)$ to $O(\mathbf{n}^4)$ constraints (depending on \mathcal{G}). The explosive growth of this relaxation scuttles most algorithms by the time $\mathbf{n} > 50$. To solve larger instances, many researchers resort to even weaker relaxations.

In contrast, we can solve the relaxation (7.6) directly using SketchyCGAL, up to $\mathbf{n} = 150$. By limiting the number of inequality constraints, via the operator \mathcal{G} , we achieve substantial reductions in resource usage. We validate our algorithm on QAPs where the exact solution is known, and we compare the performance with algorithms [107, 23] for other relaxations.

7.5.3. Rounding. Given an approximate solution to (7.6), we use a rounding method to construct a permutation. SketchyCGAL returns a matrix $\widehat{\mathbf{X}} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$ where \mathbf{U} has dimension $(\mathbf{n}^2 + 1) \times R$. We extract the first column of \mathbf{U} , discard its first entry, and reshape the remaining part into an $\mathbf{n} \times \mathbf{n}$ matrix. Then we project this matrix onto the set of permutation matrices via the Hungarian method [58, 71, 52]. This yields a feasible point $\mathbf{\Pi}$ for the problem (7.5). We repeat this procedure for all R columns of \mathbf{U} , and we pick the one that minimizes $\text{tr}(\mathbf{A}\mathbf{\Pi}\mathbf{B}\mathbf{\Pi}^*)$. This permutation gives an upper bound on the optimal value of (7.5).

7.5.4. Datasets and evaluation. We consider instances from QAPLIB [27] and TSPLIB [83] that are used in [23]. The optimal values are known, and the permutation size \mathbf{n} varies between 12 and 150. (Recall that the SDP matrix dimension $n = \mathbf{n}^2 + 1$.) We report

$$(7.7) \quad \text{relative gap \%} = \frac{\text{upper bound obtained} - \text{optimum}}{\text{optimum}} \times 100.$$

7.5.5. Solving QAPs. To solve (7.6), we cannot use the scaling (7.1) because $\|\mathcal{A}\|$ is not available; see the source code for our approach. We apply SketchyCGAL with sketch size $R = \mathbf{n}$, so the sketch uses storage $\Theta(\mathbf{n}^3)$. After rounding, a low- or medium-accuracy solution of (7.6) often provides a better permutation than a high-accuracy solution. Therefore, we applied the rounding step at iterations 2, 4, 8, 16, \dots and tracked the quality of the best permutation attained on the solution path. We stopped after (the first of) 10^6 iterations or 72 hours.

The results of this experiment appear in Figure 7.4. We compare against the best value reported by Bravo Ferreira, Khoo, and Singer in [23, Tables 4 and 6] for their CSDP method with clique size $k \in \{2, 3, 4\}$. We also include the results that [23] lists for the PATH method [107].

SketchyCGAL allows us to solve a tighter SDP relaxation of QAP than the other methods (CSDP, PATH). As a consequence, we obtain significantly smaller gaps for most instances.

8. Related work. There is a large body of literature on numerical methods for solving SDPs; see [67] for a recent survey. Here we describe the SDP literature just enough to contextualize our contributions.

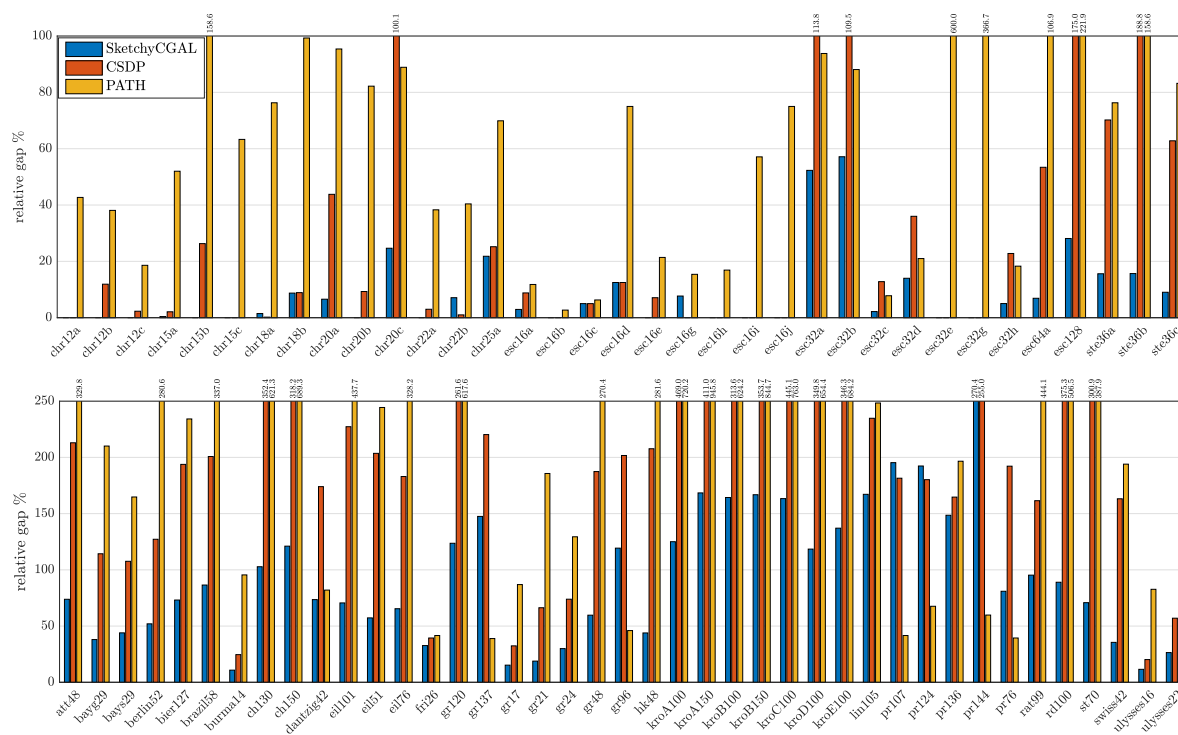


Figure 7.4. QAP relaxation: Solution quality. Using SketchyCGAL, the CSDP method [23], and the PATH method [107], we solve SDP relaxations of QAP instances from QAPLIB (top) and TSPLIB (bottom). The bars compare the cost of the computed solution against the optimal value; shorter is better. See subsection 7.5.5.

8.1. Standard methodologies. First, we outline the approaches that drive most of the reliable, general-purpose SDP software packages that are currently available.

Interior-point methods (IPMs) [75, 76, 3, 4, 53, 54] reformulate the SDP as an unconstrained problem and take an (approximate) Newton step at each iteration. In exchange, they deliver quadratic convergence. Hence IPMs are widely used to solve SDPs to high precision [56, 9, 73, 17]. Software packages include SeDuMi [90], MoSeK [70], and SDPT3 [92]. Alas, IPMs do not scale to large problems: to solve the Newton system we must store and factor large, dense matrices. A typical IPM for the SDP (2.2) requires $\Theta(n^3 + d^2n^2 + d^3)$ arithmetic operations per iteration and $\Theta(n^2 + dn + d^2)$ memory [6].

Several effective SDP solvers are based on the augmented Lagrangian paradigm [26, 98, 109, 99]. In particular, Zhao, Sun, and Toh [109] employ a semismooth Newton method and the conjugate gradient method to solve the subproblems. Their method is enhanced and implemented in the software package SDPNAL+ [100]. Augmented Lagrangian methods for SDPs typically require storage $\Omega(n^2)$.

8.2. First-order methods. First-order methods use only gradient information to solve the SDP to reduce runtime and storage requirements. We focus on *projection-free* algorithms, suitable large SDPs, that do not require full SVD computations. Major first-order methods for SDP include approaches based on the CGM [39, 63, 51] and extensions that handle more complex constraints [32, 46, 50], primal-dual subgradient algorithms [74, 105], and the matrix

multiplicative weight (MMW) method [96, 10] or, equivalently, the mirror-prox algorithm with the quantum entropy mirror map [72].

The standard CGM algorithm does not apply to the model problem (2.2) because of the affine constraint $\mathcal{A}\mathbf{X} = \mathbf{b}$. Several variants [42, 65, 87, 103] of CGM can handle affine constraints. In particular, CGAL [102] does so by applying CGM to an augmented Lagrangian formulation. We have chosen to extend CGAL because of its strong empirical performance and its robustness to inexact eigenvector computations; see [102, sect. 4-5].

Primal-dual subgradient methods perform subgradient ascent on the dual problem; the cost of each iteration is dominated by an eigenvector computation. Nesterov [74] constructs primal iterates by proximal mapping. The algorithm in Yurtsever, Tran Dinh, and Cevher [105] constructs primal iterates by an averaging technique, similar to the updates in CGM; this method involves a line search, so it requires very accurate eigenvector calculations.

The MMW method is derived by reducing the SDP to a sequence of feasibility problems. These are reformulated as a primal-dual game whose dual is an eigenvalue optimization problem. The resulting MMW algorithm can be interpreted as performing gradient descent in a dual space and using the matrix exponential map to transfer information back to the primal space. To scale this approach to larger problems, researchers have proposed linearization, random projection, sparsification techniques, and stochastic Lanczos quadrature to approximate the matrix exponential [11, 7, 80, 40, 41, 8, 13, 61, 29]. Even so, the reduction to a sequence of feasibility problems makes this technique impractical for general SDPs. We are aware of only one computational evaluation of the MMW idea [13].

8.3. Storage considerations. Almost all provably correct SDP algorithms store and operate on a full-dimensional matrix variable, so they are not suitable for very large SDPs.

Some primal-dual subgradient methods and CGM variants build an approximate solution as a convex combination of rank-one updates, so the rank of the solution does not exceed the number of iterations. This fact has led researchers to call these methods “storage-efficient,” but this claim is misleading because the algorithms require many iterations to converge.

In the conference paper [106], we observed that certain types of optimization algorithms can be combined with sketching to control storage costs. As a first example, we augmented CGM with sketching to obtain a new algorithm called SketchyCGM. This method solves a special class of low-rank matrix optimization problems that arise in statistics and machine learning. We believe that SketchyCGM is the first algorithm for this class of problems that provably succeeds with optimal storage.

To develop SketchyCGAL, we changed the base optimization algorithm (to CGAL) so that we can solve standard-form SDPs. We switched to a simpler sketching technique (the Nyström sketch) that has better empirical performance. We also analyzed how accurately to solve the eigenvalue problems to ensure that SketchyCGAL succeeds (section SM1), and we deployed an approximate eigenvalue computation method (randomized Lanczos) that meets our needs. Altogether, this effort leads to a storage-optimal algorithm that works for all standard-form SDPs with a minimum of tuning. Reducing the storage has the ancillary benefit of reducing arithmetic and communication costs, which also improves scalability.

In concurrent work with Ding [36], a subset of the authors developed a new *approximate complementarity principle* that also yields a storage-optimal algorithm for standard-form

SDPs. This approach uses a suboptimal dual point to approximate the range of the primal solution to the SDP. By compressing the primal problem to this subspace, we can solve the primal SDP with limited storage. This method, however, has more limited guarantees than SketchyCGAL. A numerical evaluation is in progress.

8.4. Nonconvex Burer–Monteiro methods. The most famous approach to low-storage semidefinite programming is the factorization heuristic proposed by Homer and Peinado [47] and refined by Burer and Monteiro (BM) [24]. The main idea is to reformulate the model problem (2.2) by expressing the psd matrix variable $\mathbf{X} = \mathbf{F}\mathbf{F}^*$ in terms of a factor $\mathbf{F} \in \mathbb{R}^{n \times R}$, where the rank parameter $R \ll n$. That is,

$$(8.1) \quad \text{minimize } \langle \mathbf{C}, \mathbf{F}\mathbf{F}^* \rangle \quad \text{subject to } \mathcal{A}(\mathbf{F}\mathbf{F}^*) = \mathbf{b}, \quad \mathbf{F}\mathbf{F}^* \in \alpha \Delta_n.$$

This approach controls storage by sacrificing convexity and the associated guarantees.

Many nonlinear programming methods have been applied to optimize (8.1). Augmented Lagrangian methods are commonly used [24, 25, 59, 85, 88]. The most popular research software based on BM factorization is Manopt [20], which implements manifold optimization algorithms including Riemannian gradient and Riemannian trust region methods [1]. Consequently, Manopt is limited to problems where the factorized formulation (8.1) defines a smooth manifold.

There has been an intense effort to establish theoretical results for the BM factorization approach. It is clear that every solution to the SDP (2.2) of rank R or less is also a solution to the factorized problem (8.1). On the other hand, (8.1) may admit local minima that are not global minima of (2.2). Some guarantees are available. For example, if \mathbf{C} is generic and the constraint set of (8.1) is a smooth manifold and $R \geq \sqrt{2(d+1)}$, then each second-order critical point of (8.1) is a global optimum [21]. A second-order critical point can be located using a Riemannian trust region method. See [19, 81, 31] for additional theoretical analysis.

The storage and arithmetic costs of solving (8.1) depend on the factorization rank R . Unfortunately, the BM method may fail when $R = o(\sqrt{d})$. Below this threshold, the BM formulation (8.1) can have spurious solutions (second-order critical points that are not globally optimal), and the bad problem instances can form a set of positive measure [97]. Hence the BM approach cannot support provably correct algorithms with storage costs better than $\Omega(n\sqrt{d})$. See subsection SM5.4 for numerical evidence.

9. Conclusion. We have presented a practical, new approach for solving SDPs at scale. Our algorithm, SketchyCGAL, combines a primal-dual optimization method with randomized linear algebra techniques to achieve unprecedented guarantees when the problem is weakly constrained and the solution is approximately low rank. We hope that our ideas lead to further algorithmic advances and support new applications of semidefinite programming.

SketchyCGAL is currently limited by the arithmetic cost of solving large eigenvalue problems to increasing accuracy. It also falters for SDPs with a large number of constraints because it depends on a primal-dual approach. Moreover, our analysis does not fully explain the observed behavior of the algorithm, including the rate of convergence of the primal variable or the convergence of the dual variable and the surrogate duality gap. These topics merit further research.

Acknowledgments. Nicolas Boumal encouraged us to reconsider the storage-optimal randomized Lanczos method (Algorithm 4.1); this algorithm simplifies the presentation while improving our theoretical and numerical results. Richard Kueng allowed us to include his results on trace normalization (step 13 in Algorithm 6.1). Irène Waldspurger provided code that generates instances of the MaxCut problem that are difficult for the Burer–Monteiro heuristic. We would also like to thank the editor and reviewers for their feedback.

REFERENCES

- [1] P.-A. ABSIL, C. BAKER, AND K. GALLIVAN, *Trust-region methods on Riemannian manifolds*, *Found. Comput. Math.*, 7 (2007), pp. 303–330.
- [2] A. Y. ALFAKIH, A. KHANDANI, AND H. WOLKOWICZ, *Solving Euclidean distance matrix completion problems via semidefinite programming*, *Comput. Optim. Appl.*, 12 (1999), pp. 13–30.
- [3] F. ALIZADEH, *Combinatorial Optimization with Interior Point Methods and Semi-Definite Matrices*, Ph.D. thesis, University of Minnesota, 1991.
- [4] F. ALIZADEH, *Interior point methods in semidefinite programming with applications to combinatorial optimization*, *SIAM J. Optim.*, 5 (1995), pp. 13–51.
- [5] F. ALIZADEH, J.-P. A. HAEBERLY, AND M. L. OVERTON, *Complementarity and nondegeneracy in semidefinite programming*, *Math. Program.*, 77 (1997), pp. 111–128.
- [6] F. ALIZADEH, J.-P. A. HAEBERLY, AND M. L. OVERTON, *Primal-dual interior point methods for semidefinite programming: Convergence rates, stability and numerical results*, *SIAM J. Optim.*, 8 (1998), pp. 746–768.
- [7] Z. ALLEN-ZHU, Y. T. LEE, AND L. ORECCHIA, *Using Optimization to Obtain a Width-Independent, Parallel, Simpler, and Faster Positive SDP Solver*, <https://arxiv.org/abs/1507.02259>, 2015.
- [8] Z. ALLEN-ZHU AND Y. LI, *Follow the Compressed Leader: Faster Online Learning of Eigenvectors and Faster MMWU*, <https://arxiv.org/abs/1701.01722>, 2017.
- [9] M. F. ANJOS AND J. B. LASSERRE, *Handbook on Semidefinite, Conic and Polynomial Optimization*, Springer, New York, 2012.
- [10] S. ARORA, E. HAZAN, AND S. KALE, *Fast algorithms for approximate semidefinite programming using the multiplicative weights update method*, in *Proceedings of the 46th Annual IEEE Symposium on Foundations of Computer Science*, Washington, DC, 2005, pp. 339–348.
- [11] S. ARORA AND S. KALE, *A combinatorial, primal-dual approach to semidefinite programs*, *J. ACM*, 63 (2016), pp. 12:1–12:35.
- [12] D. BADER, H. MEYERHENKE, P. SANDERS, AND D. WAGNER, *Graph partitioning and graph clustering*, in *10th DIMACS Implementation Challenge Workshop*, 2012, <https://www.cise.ufl.edu/research/sparse/matrices/DIMACS10/index.html>.
- [13] M. BAES, M. BÜRGISSER, AND A. NEMIROVSKI, *A randomized mirror-prox method for solving structured large-scale matrix saddle-point problems*, *SIAM J. Optim.*, 23 (2013), pp. 934–962.
- [14] R. BALAN, B. G. BODMANN, P. G. CASAZZA, AND D. EDIDIN, *Painless reconstruction from magnitudes of frame coefficients*, *J. Fourier Anal. Appl.*, 15 (2009), pp. 488–501.
- [15] R. BALAN, P. CASAZZA, AND D. EDIDIN, *On signal reconstruction without phase*, *Appl. Comput. Harmon. Anal.*, 20 (2006), pp. 345–356.
- [16] A. BARVINOK, *A Course in Convexity*, AMS, Providence, RI, 2002.
- [17] A. BEN-TAL AND A. NEMIROVSKI, *Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications*, MOS-SIAM Ser. Optim. Z, SIAM, Philadelphia, 2001.
- [18] D. P. BERTSEKAS, *Constrained Optimization and Lagrange Multiplier Methods*, *Comput. Sci. Appl. Math.*, Academic Press, New York, 1982.
- [19] S. BHOJANAPALLI, N. BOUMAL, P. JAIN, AND P. NETRAPALLI, *Smoothed analysis for low-rank solutions to semidefinite programs in quadratic penalty form*, in *Proceedings of the 31st Conference on Learning Theory*, PMLR 75, 2018, pp. 3243–3270.
- [20] N. BOUMAL, B. MISHRA, P.-A. ABSIL, AND R. SEPULCHRE, *Manopt, a MATLAB toolbox for optimization on manifolds*, *J. Mach. Learn. Res.*, 15 (2014), pp. 1455–1459.

- [21] N. BOUMAL, V. VORONINSKI, AND A. BANDEIRA, *The non-convex Burer–Monteiro approach works on smooth semidefinite programs*, in Advances in Neural Information Processing Systems 29, 2016, pp. 2757–2765.
- [22] M. BRAND, *Fast low-rank modifications of the thin singular value decomposition*, Linear Algebra Appl., 415 (2006), pp. 20 – 30.
- [23] J. F. S. BRAVO FERREIRA, Y. KHOO, AND A. SINGER, *Semidefinite programming approach for the quadratic assignment problem with a sparse graph*, Comput. Optim. Appl., 69 (2018), pp. 677–712.
- [24] S. BURER AND R. D. MONTEIRO, *A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization*, Math. Program., 95 (2003), pp. 329–357.
- [25] S. BURER AND R. D. MONTEIRO, *Local minima and convergence in low-rank semidefinite programming*, Math. Program., 103 (2005), pp. 427–444.
- [26] S. BURER AND D. VANDENBUSSCHE, *Solving lift-and-project relaxations of binary integer programs*, SIAM J. Optim., 16 (2006), pp. 726–750.
- [27] R. E. BURKARD, S. E. KARISCH, AND F. RENDL, *QAPLIB—a quadratic assignment problem library*, J. Global Optim., 10 (1997), pp. 391–403.
- [28] E. J. CANDÈS, X. LI, AND M. SOLTANOLKOTABI, *Phase retrieval from coded diffraction patterns*, Appl. Comput. Harmon. Anal., 39 (2015), pp. 277 – 299.
- [29] Y. CARMON, J. C. DUCHI, S. AARON, AND T. KEVIN, *A rank-1 sketch for matrix multiplicative weights*, in Proceedings of the 32nd Conference on Learning Theory, PMLR 99, 2019, pp. 589–623.
- [30] A. CHAI, M. MOSCOSO, AND G. PAPANICOLAOU, *Array imaging using intensity-only measurements*, Inverse Problems, 27 (2010), 015005.
- [31] D. CIFUENTES, *Burer–Monteiro Guarantees for General Semidefinite Programs*, <https://arxiv.org/abs/1904.07147>, 2019.
- [32] K. L. CLARKSON, *Coresets, sparse greedy approximation, and the Frank–Wolfe algorithm*, ACM Trans. Algorithms, 6 (2010), pp. 63:1–63:30.
- [33] J. L. CONRAD, *Marburg Virus Hemorrhagic Fever: Cytoarchitecture and Histopathology*, <https://pixnio.com/de/wissenschaft/mikroskopische-aufnahmen/hamorrhagisches-fieber-marburg-virus/cytoarchitectural-histopathologischen-erkannt-niere-probe-marburg-patient>.
- [34] C. DELORME AND S. POLJAK, *Laplacian eigenvalues and the maximum cut problem*, Math. Programming, 62 (1993), pp. 557–574.
- [35] C. DELORME AND S. POLJAK, *The performance of an eigenvalue bound on the max-cut problem in some classes of graphs*, Discrete Math., 111 (1993), pp. 145–156.
- [36] L. DING, A. YURTSEVER, V. CEVHER, J. A. TROPP, AND M. UDELL, *An Optimal-Storage Approach to Semidefinite Programming Using Approximate Complementarity*, <https://arxiv.org/abs/1902.03373>, 2019.
- [37] M. FAZEL, *Matrix Rank Minimization with Applications*, Ph.D. thesis, Stanford University, Palo Alto, CA, 2002.
- [38] J. R. FIENUP, *Phase retrieval algorithms: A comparison*, Appl. Optics, 21 (1982), pp. 2758–2769.
- [39] M. FRANK AND P. WOLFE, *An algorithm for quadratic programming*, Naval Res. Logist. Quart., 3 (1956), pp. 95–110.
- [40] D. GARBER AND E. HAZAN, *Approximating semidefinite programs in sublinear time*, in Advances in Neural Information Processing Systems 25, 2011.
- [41] D. GARBER AND E. HAZAN, *Sublinear time algorithms for approximate semidefinite programming*, Math. Programming, 158 (2016), pp. 329–361.
- [42] G. GIDEL, F. PEDREGOSA, AND S. LACOSTE-JULIEN, *Frank–Wolfe splitting via augmented Lagrangian method*, in Proceedings of the 21st International Conference on Artificial Intelligence and Statistics, PMLR 84, 2018, pp. 1456–1465.
- [43] A. GITTENS, *Topics in Randomized Numerical Linear Algebra*, Ph.D. thesis, California Institute of Technology, Pasadena, 2013.
- [44] M. X. GOEMANS AND D. P. WILLIAMSON, *Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming*, J. ACM, 42 (1995), pp. 1115–1145.
- [45] N. HALKO, P. G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions*, SIAM Rev., 53 (2011), pp. 217–288.

- [46] E. HAZAN, *Sparse approximate solutions to semidefinite programs*, in LATIN 2008: Theoretical Informatics, 2008, pp. 306–316.
- [47] S. HOMER AND M. PEINADO, *Design and performance of parallel and distributed approximation algorithms for maxcut*, J. Parallel Distrib. Comput., 46 (1997), pp. 48–61.
- [48] R. HORSTMAYER, R. Y. CHEN, X. OU, B. AMES, J. A. TROPP, AND C. YANG, *Solving ptychography with a convex relaxation*, New J. Phys., 17 (2015), 053044.
- [49] Q. HUANG, Y. CHEN, AND L. GUIBAS, *Scalable semidefinite relaxation for maximum a posterior estimation*, in Proceedings of the 31st International Conference on Machine Learning, PMLR 32, 2014, pp. 64–72.
- [50] M. JAGGI, *Revisiting Frank–Wolfe: Projection-free sparse convex optimization*, in Proceedings of the 30th International Conference on Machine Learning, PMLR 28, 2013, pp. 427–435.
- [51] L. K. JONES, *A simple lemma on greedy approximation in Hilbert space and convergence rates for projection pursuit regression and neural network training*, Ann. Statist., 20 (1992), pp. 608–613.
- [52] R. JONKER AND A. VOLGENANT, *A shortest augmenting path algorithm for dense and sparse linear assignment problems*, Computing, 38 (1987), pp. 325–340.
- [53] A. P. KAMATH AND N. K. KARMARKAR, *A continuous method for computing bounds in integer quadratic optimization problems*, J. Global Optim., 2 (1991), pp. 229–241.
- [54] A. P. KAMATH AND N. K. KARMARKAR, *An $O(nL)$ iteration algorithm for computing bounds in quadratic optimization problems*, in Complexity in Numerical Optimization, World Scientific, River Edge, NJ, 1993, pp. 254–268.
- [55] R. M. KARP, *Reducibility among combinatorial problems*, in Complexity of Computer Computations, Springer, New York, 1972, pp. 85–103.
- [56] K. KRISHNAN AND T. TERLAKY, *Interior point and semidefinite approaches in combinatorial optimization*, in Graph Theory and Combinatorial Optimization, Springer, New York, 2005.
- [57] J. KUCZYŃSKI AND H. WOŹNIAKOWSKI, *Estimating the largest eigenvalue by the power and Lanczos algorithms with a random start*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 1094–1122.
- [58] H. KUHN, *The Hungarian Method for the assignment problem*, Naval Res. Logist. Quart., 2 (1955), pp. 83–97.
- [59] B. KULIS, A. C. SURENDRAN, AND J. C. PLATT, *Fast low-rank semidefinite programming for embedding and clustering*, in Proceedings of the 11th International Conference on Artificial Intelligence and Statistics, Vol. 2, San Juan, Puerto Rico, 2007, pp. 235–242.
- [60] J. LAVAEI AND S. H. LOW, *Zero duality gap in optimal power flow problem*, IEEE Trans. Power Syst., 27 (2012), pp. 92–107.
- [61] Y. T. LEE AND S. PADMANABHAN, *An $\tilde{O}(m/\varepsilon^{3.5})$ -cost algorithm for semidefinite programs with diagonal constraints*, in Proceedings of the Conference on Learning Theory, 2020, pp. 3069–3119.
- [62] R. B. LEHOUCQ, D. C. SORENSEN, AND C. YANG, *ARPACK Users’ Guide: Solution of Large-Scale Eigenvalue Problems with Implicitly Restarted Arnoldi Methods*, Software Environ. Tools 6, SIAM, Philadelphia, 1998.
- [63] E. S. LEVITIN AND B. T. POLJAK, *Minimization methods in the presence of constraints*, Ž. Vychisl. Mat. Mat. Fiz., 6 (1966), pp. 787–823.
- [64] H. LI, G. C. LINDERMAN, A. SZLAM, K. P. STANTON, Y. KLUGER, AND M. TYGERT, *Algorithm 971: An implementation of a randomized algorithm for principal component analysis*, ACM Trans. Math. Software, 43 (2017), 28.
- [65] Y.-F. LIU, X. LIU, AND S. MA, *On the nonergodic convergence rate of an inexact augmented lagrangian framework for composite convex programming*, Math. Oper. Res., 44 (2019), pp. 632–650.
- [66] E. M. LOIOLA, N. M. M. DE ABREU, P. O. BOAVENTURA-NETTO, P. HAHN, AND T. QUERIDO, *A survey for the quadratic assignment problem*, European J. Oper. Res., 176 (2007), pp. 657–690.
- [67] A. MAJUMDAR, G. HALL, AND A. A. AHMADI, *A Survey of Recent Scalability Improvements for Semidefinite Programming with Applications to Machine Learning, Control, and Robotics*, <http://arXiv.org/abs/1908.05209>, 2019.
- [68] M. MESBAHI AND G. P. PAPAVALASSILOPOULOS, *On the rank minimization problem over a positive semidefinite linear matrix inequality*, IEEE Trans. Automat. Control, 42 (1997), pp. 239–243.
- [69] L. MIRSKY, *Symmetric gauge functions and unitarily invariant norms*, Quart. J. Math. Oxford Ser. (2), 11 (1960), pp. 50–59.

- [70] *The MOSEK Optimization Toolbox for MATLAB Manual. Version 9.0*, Mosek ApS, 2019.
- [71] J. MUNKRES, *Algorithms for the assignment and transportation problems*, J. SIAM, 5 (1957), pp. 32–38.
- [72] A. NEMIROVSKI, *Prox-method with rate of convergence $O(1/t)$ for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems*, SIAM J. Optim., 15 (2004), pp. 229–251.
- [73] Y. NESTEROV, *Introductory Lectures on Convex Optimization: A Basic Course*, Springer, New York, 2004.
- [74] Y. NESTEROV, *Primal-dual subgradient methods for convex problems*, Math. Program., 120 (2009), pp. 221–259.
- [75] Y. NESTEROV AND A. NEMIROVSKI, *Self-Concordant Functions and Polynomial Time Methods in Convex Programming*, Central Economic & Mathematic Institute report, USSR Academy of Sciences, 1990.
- [76] Y. NESTEROV AND A. NEMIROVSKI, *Interior-Point Polynomial Algorithms in Convex Programming*, SIAM, Philadelphia, 1994.
- [77] B. N. PARLETT, *The symmetric eigenvalue problem*, Classics in Appl. Math. 20, SIAM, Philadelphia, 1998, <https://doi.org/10.1137/1.9781611971163>.
- [78] P. PARRILO, *Structured Semidefinite Programs and Semialgebraic Geometry Methods in Robustness and Optimization*, Ph.D. thesis, California Institute of Technology, Pasadena, 2000.
- [79] G. PATAKI, *On the rank of extreme matrices in semidefinite programs and the multiplicity of optimal eigenvalues*, Math. Oper. Res., 23 (1998), pp. 339–358.
- [80] R. PENG AND K. TANGWONGSAN, *Faster and simpler width-independent parallel algorithms for positive semidefinite programming*, in Proceedings of the 24th Annual ACM Symposium on Parallelism in Algorithms and Architectures, 2012, pp. 101–108.
- [81] T. PUMIR, S. JELASSI, AND N. BOUMAL, *Smoothed analysis of the low-rank approach for smooth semidefinite programs*, in Advances in Neural Information Processing Systems 31, 2018, pp. 2281–2290.
- [82] A. RAGHUNATHAN, J. STEINHARDT, AND P. LIANG, *Semidefinite relaxations for certifying robustness to adversarial examples*, in Advances in Neural Information Processing Systems 31, 2018, pp. 10900–10910.
- [83] G. REINELT, *TSPLIB*, <http://comopt.ifi.uni-heidelberg.de/software/TSPLIB95/>.
- [84] D. M. ROSEN, L. CARLONE, A. S. BANDEIRA, AND J. J. LEONARD, *SE-sync: A certifiably correct algorithm for synchronization over the special euclidean group*, Int. J. Robot. Res., 38 (2019), pp. 95–125.
- [85] M. F. SAHIN, A. EFTEKHARI, A. ALACAOGLU, F. LATORRE, AND V. CEVHER, *An inexact augmented Lagrangian framework for nonconvex optimization with nonlinear constraints*, in Advances in Neural Information Processing Systems 32, 2019, pp. 13965–13977.
- [86] S. SAHNI AND T. GONZALEZ, *P-complete approximation problems*, J. ACM, 23 (1976), pp. 555–565.
- [87] A. SILVETI-FALLS, C. MOLINARI, AND J. FADILI, *Generalized Conditional Gradient with Augmented Lagrangian for Composite Minimization*, <https://arxiv.org/abs/1901.01287>, 2019.
- [88] M. SOUTO, J. D. GARCIA, AND A. VEIGA, *Exploiting Low-Rank Structure in Semidefinite Programming by Approximate Operator Splitting*, <https://arxiv.org/abs/1810.05231>, 2018.
- [89] N. SREBRO, *Learning with Matrix Factorizations*, Ph.D. thesis, Massachusetts Institute of Technology, Cambridge, 2004.
- [90] J. STURM, *Using SeDuMi 1.02, a MATLAB toolbox for optimization over symmetric cones*, Optim. Methods Softw., 11–12 (1999), pp. 625–653.
- [91] Y. SUN, Y. GUO, J. A. TROPP, AND M. UDELL, *Tensor random projection for low memory dimension reduction*, in NeurIPS Workshop on Relational Representation Learning, 2018.
- [92] K. A. TOH, M. J. TODD, AND R. H. TÜTÜNCÜ, *SDPT3—a Matlab software package for semidefinite programming, optimization methods and software*, Optim. Methods Softw., 11 (1999), pp. 545–581.
- [93] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Fixed-rank approximation of a positive-semidefinite matrix from streaming data*, in Advances in Neural Information Processing Systems 30, 2017, pp. 1225–1234.
- [94] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Practical sketching algorithms for low-rank matrix approximation*, SIAM J. Matrix Anal. Appl., 38 (2017), pp. 1454–1485.
- [95] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Streaming low-rank matrix approximation with an application to scientific simulation*, SIAM J. Sci. Comput., 41 (2019), pp. A2430–A2463.

- [96] K. TSUDA, G. RÄTSCH, AND M. K. WARMUTH, *Matrix exponentiated gradient updates for on-line learning and Bregman projections*, *J. Mach. Learn. Res.*, 6 (2005), pp. 995–1018.
- [97] I. WALDSPURGER AND A. WATERS, *Rank optimality for the Burer–Monteiro factorization*, <https://arxiv.org/abs/1812.03046>, 2018.
- [98] Z. WEN, D. GOLDFARB, S. MA, AND K. SCHEINBERG, *Row by Row Methods for Semidefinite Programming*, IEOB report, Columbia University, 2009.
- [99] Z. WEN, D. GOLDFARB, AND W. YIN, *Alternating direction augmented Lagrangian methods for semidefinite programming*, *Math. Program. Comput.*, 2 (2010), pp. 203–230.
- [100] L. YANG, D. SUN, AND K.-C. TOH, *SDPNAL+: A majorized semismooth Newton-CG augmented Lagrangian method for semidefinite programming with nonnegative constraints*, *Math. Program. Comput.*, 7 (2015), pp. 331–366.
- [101] Y. YE, *GSet Random Graphs*, <https://www.cise.ufl.edu/research/sparse/matrices/Gset/>.
- [102] A. YURTSEVER, O. FERCOQ, AND V. CEVHER, *A conditional-gradient-based augmented Lagrangian framework*, in *Proceedings of the 36th International Conference on Machine Learning*, PMLR 97, 2019, pp. 7272–7281.
- [103] A. YURTSEVER, O. FERCOQ, F. LOCATELLO, AND V. CEVHER, *A conditional gradient framework for composite convex minimization with applications to semidefinite programming*, in *Proceedings of the 35th International Conference on Machine Learning*, PMLR 80, 2018, pp. 5727–5736.
- [104] A. YURTSEVER, Y.-P. HSIEH, AND V. CEVHER, *Scalable convex methods for phase retrieval*, in *Proceedings of the 6th International IEEE Workshop on Computational Advances in Multi-Sensor Adaptive Processing*, 2015, pp. 381–384.
- [105] A. YURTSEVER, Q. TRAN DINH, AND V. CEVHER, *A universal primal-dual convex optimization framework*, in *Advances in Neural Information Processing Systems 28*, 2015, pp. 3150–3158.
- [106] A. YURTSEVER, M. UDELL, J. TROPP, AND V. CEVHER, *Sketchy decisions: Convex low-rank matrix optimization with optimal storage*, in *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, PMLR, 54, 2017, pp. 1188–1196.
- [107] M. ZASLAVSKIY, F. BACH, AND J. VERT, *A path following algorithm for the graph matching problem*, *IEEE Trans. Pattern Anal. Mach. Intel.*, 31 (2009), pp. 2227–2242.
- [108] Q. ZHAO, S. E. KARISCH, F. RENDL, AND H. WOLKOWICZ, *Semidefinite programming relaxations for the quadratic assignment problem*, *J. Combin. Optim.*, 2 (1998), pp. 71–109.
- [109] X.-Y. ZHAO, D. SUN, AND K.-C. TOH, *A Newton-CG augmented Lagrangian method for semidefinite programming*, *SIAM J. Optim.*, 20 (2010), pp. 1737–1765.
- [110] G. ZHENG, R. HORSTMAYER, AND C. YANG, *Wide-field, high-resolution Fourier ptychography microscopy*, *Nat. Photonics*, 7 (2013), pp. 739–745.