COLUMN SUBSET SELECTION, MATRIX FACTORIZATION, AND EIGENVALUE OPTIMIZATION

JOEL A. TROPP

1891

Technical Report No. 2008-02 July 2008

COLUMN SUBSET SELECTION, MATRIX FACTORIZATION, AND EIGENVALUE OPTIMIZATION

J. A. TROPP

ABSTRACT. Given a fixed matrix, the problem of column subset selection requests a column submatrix that has favorable spectral properties. Most research from the algorithms and numerical linear algebra communities focuses on a variant called rank-revealing QR, which seeks a well-conditioned collection of columns that spans the (numerical) range of the matrix. The functional analysis literature contains another strand of work on column selection whose algorithmic implications have not been explored. In particular, a celebrated result of Bourgain and Tzafriri demonstrates that each matrix with normalized columns contains a large column submatrix that is exceptionally well conditioned. Unfortunately, standard proofs of this result cannot be regarded as algorithmic.

This paper presents a randomized, polynomial-time algorithm that produces the submatrix promised by Bourgain and Tzafriri. The method involves random sampling of columns, followed by a matrix factorization that exposes the well-conditioned subset of columns. This factorization, which is due to Grothendieck, is regarded as a central tool in modern functional analysis. The primary novelty in this work is an algorithm, based on eigenvalue minimization, for constructing the Grothendieck factorization. These ideas also result in a novel approximation algorithm for the $(\infty,1)$ norm of a matrix, which is generally NP-hard to compute exactly. As an added bonus, this work reveals a surprising connection between matrix factorization and the famous MAXCUT semidefinite program.

Date: 26 June 2008.

JAT is with Applied and Computational Mathematics, MC 217-50, California Inst. Technology, Pasadena, CA 91125-5000. E-mail: jtropp@acm.caltech.edu. Supported in part by ONR award no. N000140810883.

1. Introduction

Column subset selection refers to the challenge of extracting from a matrix a column submatrix that has some distinguished property. These properties commonly involve conditions on the spectrum of the submatrix. The most familiar example is probably rank-revealing QR, which seeks a well-conditioned collection of columns that spans the (numerical) range of the matrix [GE96].

The literature on geometric functional analysis contains several fundamental theorems on column subset selection that have not been discussed by the algorithms community or the numerical linear algebra community. These results are phrased in terms of the *stable rank* of a matrix:

$$\operatorname{st.rank}(\boldsymbol{A}) = \frac{\|\boldsymbol{A}\|_{\operatorname{F}}^2}{\|\boldsymbol{A}\|^2}$$

where $\|\cdot\|_{F}$ is the Frobenius norm and $\|\cdot\|$ is the spectral norm. The stable rank can be viewed as an analytic surrogate for the algebraic rank. Indeed, express the two norms in terms of singular values to obtain the relation

st.
$$rank(\mathbf{A}) \leq rank(\mathbf{A})$$
.

In this bound, equality occurs (for example) when the columns of A are identical or when the columns of A are orthonormal. As we will see, the stable rank is tightly connected with the number of (strongly) linearly independent columns we can extract from a matrix.

Before we continue, let us instate some notation. We say that a matrix is *standardized* when its columns have unit ℓ_2 norm. The *j*th column of a matrix \boldsymbol{A} is denoted by \boldsymbol{a}_j . For a subset τ of column indices, we write \boldsymbol{A}_{τ} for the column submatrix indexed by τ . Likewise, given a square matrix \boldsymbol{H} , the notation $\boldsymbol{H}_{\tau \times \tau}$ refers to the principal submatrix whose rows and columns are listed in τ . The pseudoinverse \boldsymbol{D}^{\dagger} of a diagonal matrix \boldsymbol{D} is formed by reciprocating the nonzero entries. As usual, we write $\|\cdot\|_p$ for the ℓ_p vector norm. The *condition number* of a matrix is the quantity

$$\kappa(oldsymbol{A}) = \max \left\{ rac{\|oldsymbol{A}oldsymbol{x}\|_2}{\|oldsymbol{A}oldsymbol{y}\|_2} : \|oldsymbol{x}\|_2 = \|oldsymbol{y}\|_2 = 1
ight\}.$$

Finally, upright letters (c, C, K, ...) refer to positive, universal constants that may change from appearance to appearance.

The first theorem, due to Kashin and Tzafriri, shows that each matrix with standardized columns contains a large column submatrix that has small spectral norm [Ver01, Thm. 2.5].

Theorem 1.1 (Kashin–Tzafriri). Suppose A is standardized. Then there is a set τ of column indices for which

$$|\tau| \ge \text{st. rank}(\boldsymbol{A})$$
 and $\|\boldsymbol{A}_{\tau}\| \le C$.

In fact, much more is true. Combining Theorem 1.1 with the celebrated restricted invertibility result of Bourgain and Tzafriri [BT87, Thm. 1.2], we find that every standardized matrix contains a large column submatrix whose *condition number* is small.

Theorem 1.2 (Bougain–Tzafriri). Suppose A is standardized. Then there is a set τ of column indices for which

$$|\tau| \ge c \cdot \text{st. rank}(\boldsymbol{A}) \quad and \quad \kappa(\boldsymbol{A}_{\tau}) \le \sqrt{3}.$$

Theorem 1.2 yields the best general result [BT91, Thm. 1.1] on the Kadison–Singer conjecture, the major open question in operator theory. To display its strength, let us consider two extreme examples.

- (1) When \boldsymbol{A} has identical columns, every collection of two or more columns is singular. Theorem 1.2 guarantees a well-conditioned submatrix \boldsymbol{A}_{τ} with $|\tau| = 1$, which is optimal.
- (2) When A has n orthonormal columns, the full matrix is perfectly conditioned. Theorem 1.2 guarantees a well-conditioned submatrix A_{τ} with $|\tau| \geq cn$, which lies within a constant factor of optimal.

The stable rank allows Theorem 1.2 to interpolate between the two extremes. Subsequent research established that the stable rank is intrinsic to the problem of finding well-conditioned submatrices. We postpone a more detailed discussion of this point until Section 6.

1.1. Contributions. Although Theorems 1.1 and 1.2 would be very useful in computational applications, we cannot regard current proofs as constructive. The goal of this paper is to establish the following novel, algorithmic claim.

Theorem 1.3. There are randomized, polynomial-time algorithms for producing the sets guaranteed by Theorem 1.1 and by Theorem 1.2.

This result is significant because no known algorithm for column subset selection is guaranteed to produce a submatrix whose condition number has constant order. See [BDM08] for a recent overview of that literature. The present work has other ramifications with independent interest.

- We develop algorithms for computing the matrix factorizations of Pietsch and Grothendieck, which are regarded as basic instruments in modern functional analysis [Pis86].
- The methods for computing these factorizations lead to new approximation algorithms for two NP-hard matrix norms. (See Remarks 3.2 and 5.6.)
- We identify an intriguing connection between Pietsch factorization and the MAXCUT semi-definite program [GW95].
- 1.2. **Overview.** We focus on the algorithmic version of the Kashin–Tzafriri theorem because it highlights all the essential concepts while minimizing irrelevant details. Section 2 outlines a proof of this result, emphasizing where new algorithmic machinery is required. The missing link turns out to be a computational method for producing a certain matrix factorization. Section 3 reformulates the factorization problem as an eigenvalue minimization, which can be completed with standard techniques. In Section 4, we exhibit a randomized algorithm that delivers the submatrix promised by Kashin–Tzafriri. In Section 5, we traverse a similar route to develop an algorithmic version of Bourgain–Tzafriri. Section 6 provides more details about the stable rank and describes directions for future work. Appendix A contains some key estimates on the norms of random submatrices, and Appendix B outlines a simple computational procedure for solving the eigenvalue optimization problems that arise in our work.

2. The Kashin-Tzafriri Theorem

The proof of the Kashin–Tzafriri theorem proceeds in two steps. First, we select a random set of columns with appropriate cardinality. Second, we use a matrix factorization to identify and remove redundant columns that inflate the spectral norm. The proof gives strong hints about how a computational procedure might work, even though it is not constructive.

2.1. **Intuitions.** We would like to think that a random submatrix inherits its share of the norm of the entire matrix. In other words, if we were to select a tenth of the columns, we might hope to reduce the norm by a factor of ten. Unfortunately, this intuition is meretricious.

Indeed, random selection does not necessarily reduce the spectral norm at all. The essential reason emerges when we consider the "double identity," the $m \times 2m$ matrix $\mathbf{A} = \begin{bmatrix} \mathbf{I} \mid \mathbf{I} \end{bmatrix}$. Suppose we draw s random columns from \mathbf{A} without replacement. The probability that all s columns are distinct is

$$\frac{2m-2}{2m-1} \times \frac{2m-4}{2m-2} \times \dots \times \frac{2m-2(s-1)}{2m-(s-1)} \le \prod_{j=0}^{s-1} \left(1 - \frac{j}{2m}\right) \approx \exp\left\{-\sum_{j=0}^{s-1} \frac{j}{2m}\right\} \approx e^{-s^2/4m}.$$

Therefore, when $s = \Omega(\sqrt{m})$, sampling almost always produces a submatrix with at least one duplicated column. A duplicated column means that the norm of the submatrix is $\sqrt{2}$, which equals the norm of the full matrix, so no reduction takes place.

Nevertheless, a randomly chosen set of columns from a standardized matrix typically *contains* a large set of columns that has small norm. We will see that the desired subset is exposed by factoring the random submatrix. This factorization, which was invented by Pietsch, is regarded as a basic instrument in modern functional analysis.

2.2. The $(\infty, 2)$ operator norm. Although sampling does not necessarily reduce the spectral norm, it often reduces other matrix norms. Define the natural norm on linear operators from ℓ_{∞} to ℓ_2 via the expression

$$\|{\pmb B}\|_{\infty \to 2} = \max\{\|{\pmb B}{\pmb x}\|_2 : \|{\pmb x}\|_{\infty} = 1\}.$$

An immediate consequence is that $\|\boldsymbol{B}\|_{\infty\to 2} \leq \sqrt{s} \|\boldsymbol{B}\|$ for each matrix \boldsymbol{B} with s columns. Equality can obtain in this bound.

The exact calculation of the $(\infty, 2)$ operator norm is computationally difficult. Results of Rohn [Roh00] imply that there is a class of positive semidefinite matrices for which it is NP-hard to estimate $\|\cdot\|_{\infty\to 2}$ within an absolute tolerance. Nevertheless, we will see that the norm can be approximated in polynomial time up to a small relative error. (See Remark 3.2.)

As we have intimated, the $(\infty, 2)$ norm can often be reduced by random selection. The following theorem requires some heavy lifting, which we delegate to Appendix A.2.

Theorem 2.1. Suppose A is a standardized matrix with n columns. Choose

$$s \leq \lceil 2 \operatorname{st.rank}(\boldsymbol{A}) \rceil$$
,

and draw a uniformly random subset σ with cardinality s from $\{1, 2, ..., n\}$. Then

$$\mathbb{E} \|\boldsymbol{A}_{\sigma}\|_{\infty \to 2} \le 7\sqrt{s}.$$

In particular, $\|\mathbf{A}_{\sigma}\|_{\infty \to 2} \leq 8\sqrt{s}$ with probability at least 1/8.

2.3. **Pietsch factorization.** We cannot exploit the bound in Theorem 2.1 unless we have a way to connect the $(\infty, 2)$ norm with the spectral norm. To that end, let us recall one of the landmark theorems of functional analysis.

Theorem 2.2 (Pietsch Factorization). Each matrix B can be factored as B = TD where

- **D** is a nonnegative, diagonal matrix with trace(D^2) = 1, and
- $\bullet \|B\|_{\infty \to 2} \le \|T\| \le K_{\mathbf{P}} \|B\|_{\infty \to 2}.$

This result follows from the little Grothendieck theorem [Pis86, Sec. 5b] and the Pietsch factorization theorem [Pis86, Cor. 1.8]. The standard proof produces the factorization using an abstract separation argument that offers no algorithmic insight. The value of the constant is available.

- When the scalar field is real, we have $K_P(\mathbb{R}) = \sqrt{\pi/2} \approx 1.25$.
- When the scalar field is complex, we have $K_P(\mathbb{C}) = \sqrt{4/\pi} \approx 1.13$.

A major application of Pietsch factorization is to identify a submatrix with controlled spectral norm. The following proposition describes the procedure.

Proposition 2.3. Suppose B is a matrix with s columns. Then there is a set τ of column indices for which

$$| au| \geq rac{s}{2} \quad and \quad \|oldsymbol{B}_{ au}\| \leq \mathrm{K}_{\mathrm{P}} \sqrt{rac{2}{s}} \, \|oldsymbol{B}\|_{\infty o 2} \, .$$

Proof. Consider a Pietsch factorization B = TD, and define

$$\tau = \{j : d_{jj}^2 \le 2/s\}.$$

Since $\sum d_{ij}^2 = 1$, Markov's inequality implies that $|\tau| \geq s/2$. We may calculate that

$$||B_{\tau}|| = ||TD_{\tau}|| \le ||T|| \cdot ||D_{\tau}|| \le K_{\mathrm{P}} ||B||_{\infty \to 2} \cdot \sqrt{2/s}.$$

This completes the proof.

2.4. **Proof of Kashin–Tzafriri.** With these results at hand, we easily complete the proof of the Kashin–Tzafriri theorem. Suppose \mathbf{A} is a standardized matrix with n columns. Assume that st. rank(\mathbf{A}) $\leq n/2$. Otherwise, the spectral norm $\|\mathbf{A}\| \leq \sqrt{2}$, so we may select $\tau = \{1, 2, ..., n\}$.

According to Theorem 2.1, there is a subset σ of column indices for which

$$|\sigma| \ge 2 \text{ st. rank}(\boldsymbol{A}) \quad \text{and} \quad \|\boldsymbol{A}_{\sigma}\|_{\infty \to 2} \le 8\sqrt{|\sigma|}.$$

Apply Proposition 2.3 to the matrix $B = A_{\sigma}$ to obtain a subset τ inside σ for which

$$|\tau| \geq \frac{|\sigma|}{2}$$
 and $\|\boldsymbol{B}_{\tau}\| \leq K_{P} \sqrt{\frac{2}{|\sigma|}} \|\boldsymbol{B}\|_{\infty \to 2}$.

Since $B_{\tau} = A_{\tau}$ and $K_{P} \leq \sqrt{\pi/2}$, these bounds reveal the advertised conclusion:

$$|\tau| \ge \text{st.rank}(\boldsymbol{A})$$
 and $\|\boldsymbol{A}_{\tau}\| < 15$.

At this point, we take a step back and notice that this proof is nearly algorithmic. It is straightforward to perform the random selection described in Theorem 2.1. Provided that we know a Pietsch factorization of the matrix \boldsymbol{B} , we can easily carry out the column selection of Proposition 2.3. Therefore, we need only develop an algorithm for computing the Pietsch factorization to reach an effective version of the Kashin–Tzafriri theorem.

3. Pietsch Factorization via Convex Optimization

The main novelty is to demonstrate that we can produce a Pietsch factorization by solving a convex programming problem. Remarkably, the resulting optimization is the dual of the famous MAXCUT semidefinite program [GW95], for which many polynomial-time algorithms are available.

3.1. Pietsch and eigenvalues. The next theorem, which serves as the basis for our computational method, demonstrates that Pietsch factorizations have an intimate relationship with the eigenvalues of a related matrix. In the sequel, we reserve the letter D for a nonnegative, diagonal matrix with trace(D^2) = 1, and we write λ_{max} for the algebraically maximal eigenvalue of a Hermitian matrix.

Theorem 3.1. The factorization B = TD satisfies $||T|| \le \alpha$ if and only if D satisfies

$$\lambda_{\max}(\boldsymbol{B}^*\boldsymbol{B} - \alpha^2 \boldsymbol{D}^2) \le 0.$$

In particular, if no **D** verifies this bound, then no factorization B = TD admits $||T|| \leq \alpha$.

Proof. Assume B has a factorization B = TD with $||T|| \leq \alpha$. We have the chain of implications

$$\begin{aligned} \boldsymbol{B} &= \boldsymbol{T} \boldsymbol{D} &\implies & \|\boldsymbol{B} \boldsymbol{x}\|_2^2 = \|\boldsymbol{T} \boldsymbol{D} \boldsymbol{x}\|_2^2 \\ &\implies & \|\boldsymbol{B} \boldsymbol{x}\|_2^2 \leq \alpha^2 \|\boldsymbol{D} \boldsymbol{x}\|_2^2 & \forall \boldsymbol{x} \\ &\implies & \boldsymbol{x}^* \boldsymbol{B}^* \boldsymbol{B} \boldsymbol{x} \leq \alpha^2 \boldsymbol{x}^* \boldsymbol{D}^2 \boldsymbol{x} & \forall \boldsymbol{x} \\ &\implies & \boldsymbol{x}^* (\boldsymbol{B}^* \boldsymbol{B} - \alpha^2 \boldsymbol{D}^2) \boldsymbol{x} \leq 0 & \forall \boldsymbol{x} \\ &\implies & \boldsymbol{B}^* \boldsymbol{B} - \alpha^2 \boldsymbol{D}^2 \leqslant \mathbf{0}, \end{aligned}$$

where \leq denotes the semidefinite, or Löwner, ordering on Hermitian matrices.

Conversely, assume we are provided the inequality

$$B^*B - \alpha^2 D^2 \leq 0. \tag{3.1}$$

First, we claim that any zero entry in D corresponds with a zero column of B. To check this point, suppose that $d_{ij} = 0$ for an index j. The relation (3.1) requires that

$$0 \ge (\boldsymbol{B}^* \boldsymbol{B} - \alpha^2 \boldsymbol{D}^2)_{jj} = \boldsymbol{b}_j^* \boldsymbol{b}_j.$$

This inequality is impossible unless $b_j = 0$. To continue, set $T = BD^{\dagger}$, and observe that B = TD because the zero entries of D correspond with zero columns of B. Therefore, we may factor the diagonal matrix out from (3.1) to reach

$$D(T^*T - \alpha^2 P)D \leq 0.$$

where the matrix $P = DD^{\dagger}$ is an orthogonal projector. Sylvester's theorem on inertia [HJ85, Thm. 4.5.8] ensures that $T^*T - \alpha^2 P \leq 0$. Since P is a projector, this relation implies that

$$T^*T \preceq \alpha^2 P \preceq \alpha^2 I$$
.

We conclude that $||T|| \leq \alpha$.

3.2. Factorization via optimization. Recall that the maximum eigenvalue is a convex function on the space of Hermitian matrices, so it can be minimized in polynomial time [LO96]. We are led to consider the convex program

min
$$\lambda_{\max}(\mathbf{B}^*\mathbf{B} - \alpha^2\mathbf{F})$$
 subject to trace(\mathbf{F}) = 1, \mathbf{F} diagonal, and $\mathbf{F} \ge \mathbf{0}$. (3.2)

Owing to Theorem 3.1, there exists a factorization B = TD with $||T|| \le \alpha$ if and only if the value of (3.2) is nonpositive.

Now, if F is a feasible point of (3.2) with a nonpositive objective value, we can factorize

$$B = TD$$
 with $D = F^{1/2}$, $T = BD^{\dagger}$, and $||T|| \le \alpha$.

In fact, it is not necessary to solve (3.2) to optimality. Suppose B has s columns, and assume we have identified a feasible point F with a (positive) objective value η . That is,

$$\lambda_{\max}(\boldsymbol{B}^*\boldsymbol{B} - \alpha^2 \boldsymbol{F}) \le \eta.$$

Rearranging this relation, we reach

$$\lambda_{\max} \left[\mathbf{B}^* \mathbf{B} - (\alpha^2 + \eta s) \widetilde{\mathbf{F}} \right] \le 0 \text{ where } \widetilde{\mathbf{F}} = \frac{1}{\alpha^2 + \eta s} (\alpha^2 \mathbf{F} + \eta \mathbf{I}).$$

Since \widetilde{F} is positive and diagonal with trace $(\widetilde{F}) = 1$, we obtain the factorization

$$B = TD$$
 with $D = \widetilde{F}^{1/2}$, $T = BD^{-1}$, and $||T|| < \sqrt{\alpha^2 + \eta s}$.

To select a target value for the parameter α , we look to the proof of the Kashin–Tzafriri theorem. If \boldsymbol{B} has s columns, then $\alpha = 8 K_P \sqrt{s}$ is an appropriate choice. Furthermore, since the argument only uses the bound $\|\boldsymbol{T}\| = O(\sqrt{s})$, it suffices to solve (3.2) with precision $\eta = O(1)$.

3.3. Other formulations. In a general setting, a target value for α is not likely to be available. Let us exhibit an alternative formulation of (3.2) that avoids this inconvenience.

min
$$\lambda_{\max}(B^*B - E) + \operatorname{trace}(E)$$
 subject to E diagonal, $E \ge 0$. (3.3)

Suppose α_{\star} is the minimal value of ||T|| achievable in any Pietsch factorization B = TD. It can be shown that α_{\star}^2 is the value of (3.3) and that each optimizer E_{\star} satisfies trace(E_{\star}) = α_{\star}^2 . As such, we can construct an optimal Pietsch factorization from a minimizer:

$$m{B} = m{T}m{D} \quad ext{with} \quad m{D} = (m{E}_{\star}/\operatorname{trace}(m{E}_{\star}))^{1/2}, \quad m{T} = m{B}m{D}^{\dagger}, \quad ext{and} \quad \|m{T}\| = lpha_{\star}.$$

The dual of (3.3) is the semidefinite program

$$\max \langle B^*B, Z \rangle$$
 subject to $\operatorname{diag}(Z) = \mathbf{I} \text{ and } Z \geq 0.$ (3.4)

This is the famous MAXCUT semidefinite program [GW95]. We find an unexpected connection between Pietsch factorization and the problem of partitioning nodes of a graph.

Given a dual optimum, we can easily construct a primal optimum by means of the complementary slackness condition [Ali95, Thm. 2.10]. Indeed, each feasible optimal pair (E_{\star}, Z_{\star}) satisfies $Z_{\star}(B^*B - E_{\star}) = 0$. Examining the diagonal elements of this matrix equation, we find that

$$E_{\star} = \operatorname{diag}(E_{\star}) = \operatorname{diag}(ZE_{\star}) = \operatorname{diag}(Z_{\star}B^{*}B)$$

owing to the constraint diag(\mathbf{Z}_{\star}) = **I**. Obtaining a dual optimum from a primal optimum, however, requires more ingenuity.

Remark 3.2. According to Theorem 2.2 and the discussion here, the optimal value of (3.3) overestimates $\|B\|_{\infty\to 2}^2$ by a multiplicative factor no greater than K_P^2 . As a result, the optimization problem (3.3) can be used to design an approximation algorithm for $(\infty, 2)$ norms.

3.4. Algorithmic aspects. The purpose of this paper is not to rehash methods for solving a standard optimization problem, so we keep this discussion brief. It is easy to see that (3.2) can be framed as a (nonsmooth) convex optimization over the probability simplex. Appendix B outlines an elegant technique, called Entropic Mirror Descent [BT03], designed specifically for this class of problems. Although the EMD algorithm is (theoretically) not the most efficient approach to (3.2), preliminary experiments suggest that its empirical performance rivals more sophisticated techniques.

For a concrete time bound, we refer to Alizadeh's work on primal-dual potential reduction methods for semidefinite programming [Ali95]. When \mathbf{B} has dimension $m \times s$, the cost of forming $\mathbf{B}^*\mathbf{B}$ is at most $\mathrm{O}(s^2m)$. Then the cost of solving (3.4) is no more than $\mathrm{O}(s^{3.5})$, where the tilde indicates that log-like factors are suppressed.

4. An Algorithm for Kashin-Tzafriri

At this point, we have a massed the matériel necessary to deploy an algorithm that constructs the set τ promised by the Kashin–Tzafriri theorem. The procedure appears on page 11 as Algorithm 1. The following result describes its performance.

Theorem 4.1. Suppose A is an $m \times n$ standardized matrix. With probability at least 4/5, Algorithm 1 produces a set $\tau = \tau_{\star}$ of column indices for which

$$|\tau| \ge \frac{1}{2} \operatorname{st.rank}(\boldsymbol{A}) \quad and \quad \|\boldsymbol{A}_{\tau}\| \le 15.$$

The computational cost is bounded by $\widetilde{O}(|\tau|^2 m + |\tau|^{3.5})$.

Remarkably, Algorithm 1 is sublinear in the size of the matrix when st.rank(\mathbf{A}) = o($n^{1/3.5}$). Better methods for solving (3.2) would strengthen this bound.

Proof. According to Section 2, the procedure NORM-REDUCE has failure probability less than 7/8 when $s \leq 2$ st. rank(\mathbf{A}). The probability the inner loop fails to produce an acceptable set τ_{\star} of size s/2 is at most $(7/8)^{8\log_2(s)}$. So the probability the algorithm fails before $s \geq \text{st. rank}(\mathbf{A})$ is at most

$$\sum_{j=2}^{\infty} (7/8)^{8j} = \frac{(7/8)^{16}}{1 - (7/8)^8} < 0.2.$$

With constant probability, we obtain a set τ_{\star} with cardinality at least st.rank(\mathbf{A})/2.

The cost of the procedure NORM-REDUCE is dominated by the cost of the Pietsch factorization, which is $\widetilde{O}(s^2m + s^{3.5})$ for a fixed s. Summing over s and k, we find that the total cost of all the invocations of NORM-REDUCE is dominated (up to logarithmic factors) by the cost of the final invocation, during which the parameter $s < 2 |\tau_{\star}|$.

An estimate of the spectral norm of A_{τ} can be obtained as a by-product of solving (3.2). Indeed, Proposition 2.3 and the discussion in Section 3.2 show that we can bound the spectral norm in terms of the parameter α and the objective value obtained in (3.2).

5. The Bourgain-Tzafriri Theorem

Our proof of the Bourgain-Tzafriri theorem is almost identical in structure with the proof of the Kashin-Tzafriri theorem. This streamlined argument appears to be simpler than all previously published approaches, but it contains no significant conceptual innovations. Our discussion culminates in an algorithm remarkably similar to Algorithm 1.

5.1. **Preliminary results.** Suppose A is a standardized matrix with n columns. We will work instead with a related matrix $H = A^*A - I$, which is called the *hollow Gram matrix*. The advantage of considering the hollow Gram matrix is that we can perform column selection on A simply by reducing the norm of H.

Proposition 5.1. Suppose \mathbf{A} is a standardized matrix with hollow Gram matrix \mathbf{H} . If τ is a set of column indices for which $\|\mathbf{H}_{\tau \times \tau}\| \leq 0.5$, then $\kappa(\mathbf{A}_{\tau}) \leq \sqrt{3}$.

Proof. The hypothesis $\|\boldsymbol{H}_{\tau \times \tau}\| \leq 0.5$ implies that the eigenvalues of $\boldsymbol{H}_{\tau \times \tau}$ lie in the range [-0.5, 0.5]. Since $\boldsymbol{H}_{\tau \times \tau} = \boldsymbol{A}_{\tau}^* \boldsymbol{A}_{\tau} - \mathbf{I}$, the eigenvalues of $\boldsymbol{A}_{\tau}^* \boldsymbol{A}_{\tau}$ fall in the interval [0.5, 1.5]. An equivalent condition is that $0.5 \leq \|\boldsymbol{A}_{\tau} \boldsymbol{x}\|_2^2 \leq 1.5$ whenever $\|\boldsymbol{x}\|_2 = 1$. We conclude that

$$\kappa(\boldsymbol{A}_{\tau}) = \max \left\{ \frac{\|\boldsymbol{A}_{\tau}\boldsymbol{x}\|_{2}}{\|\boldsymbol{A}_{\tau}\boldsymbol{y}\|_{2}} : \|\boldsymbol{x}\|_{2} = \|\boldsymbol{y}\|_{2} = 1 \right\} \leq \sqrt{\frac{1.5}{0.5}} = \sqrt{3}.$$

Thus, a norm bound for $H_{\tau \times \tau}$ yields a condition number bound for A_{τ} .

As we mentioned before, random selection may reduce other norms even if it does not reduce the spectral norm. Define the natural norm on linear maps from ℓ_{∞} to ℓ_1 by the formula

$$\|\boldsymbol{G}\|_{\infty \to 1} = \max\{\|\boldsymbol{G}\boldsymbol{x}\|_1 : \|\boldsymbol{x}\|_{\infty} = 1\}.$$

This norm is closely related to the cut norm, which plays a starring role in graph theory [AN04]. For a general $s \times s$ matrix G, the best inequality between the $(\infty, 1)$ norm and the spectral norm is $\|G\|_{\infty \to 1} \le s \|G\|$. Rohn [Roh00] has established that there is a class of positive semidefinite, integer matrices for which it is NP-hard to determine the $(\infty, 1)$ norm within an absolute tolerance of 1/2. Nevertheless, it can be approximated within a small relative factor in polynomial time [AN04].

The $(\infty, 1)$ norm decreases when we randomly sample a principal submatrix. The following result, which we establish in Appendix A.4, is a direct consequence of Rudelson and Vershynin's work on the cut norm of random submatrices [RV07, Thm. 1.5].

Theorem 5.2. Suppose A is an n-column standardized matrix with hollow Gram matrix H. Choose

$$s \leq [\mathbf{c} \cdot \operatorname{st.rank}(\mathbf{A})],$$

and draw a uniformly random subset σ with cardinality s from $\{1, 2, ..., n\}$. Then

$$\mathbb{E} \| \boldsymbol{H}_{\sigma \times \sigma} \|_{\infty \to 1} \le \frac{s}{9}.$$

In particular, $\|\mathbf{H}_{\sigma \times \sigma}\|_{\infty \to 1} \le s/8$ with probability at least 1/9.

To connect the $(\infty, 1)$ norm with the spectral norm, we call on the celebrated factorization of Grothendieck [Pis86, p. 56].

Theorem 5.3 (Grothendieck Factorization). Each matrix G can be factored as $G = D_1TD_2$ where

- (1) D_i is a nonnegative, diagonal matrix with trace $(D_i^2) = 1$ for i = 1, 2, and
- (2) $\|G\|_{\infty \to 1} \le \|T\| \le K_G \|G\|_{\infty \to 1}$.

When G is Hermitian, we may take $D_1 = D_2$.

The precise value of the Grothendieck constant K_G remains an outstanding open question, but it is known to depend on the scalar field [Pis86, Sec. 5e].

- When the scalar field is real, $1.570 \le \pi/2 \le K_G(\mathbb{R}) \le \pi/(2\log(1+\sqrt{2})) \le 1.783$.
- When the scalar field is complex, $1.338 \leq K_G(\mathbb{C}) \leq 1.405$.

For positive semidefinite G, the real (resp., complex) Grothendieck constant equals the square of the real (resp., complex) Pietsch constant because $\|\boldsymbol{B}^*\boldsymbol{B}\|_{\infty\to 1} = \|\boldsymbol{B}\|_{\infty\to 2}^2$. The following proposition describes the role of the Grothendieck factorization in the selection of

submatrices with controlled spectral norm.

Proposition 5.4. Suppose G is an $s \times s$ Hermitian matrix. There is a set τ of column indices for which

$$| au| \geq rac{s}{2} \quad and \quad \|oldsymbol{G}_{ au imes au}\| \leq rac{2 K_G}{s} \, \|oldsymbol{G}\|_{\infty o 1} \, .$$

Proof. Consider a Grothendieck factorization G = DTD, and identify $\tau = \{j : d_{ij}^2 \le s/2\}$. The remaining details echo the proof of Proposition 2.3.

5.2. **Proof of Bourgain-Tzafriri.** Suppose A is a standardized matrix with n columns, and consider its hollow Gram matrix H. Theorem 5.2 provides a set σ for which

$$|\sigma| \ge c \cdot \text{st. rank}(\boldsymbol{A}) \quad \text{and} \quad \|\boldsymbol{H}_{\sigma \times \sigma}\|_{\infty \to 1} \le \frac{s}{8}.$$

Apply Proposition 5.4 to the $s \times s$ matrix $G = H_{\sigma \times \sigma}$ to obtain a further subset τ inside σ with

$$| au| \geq rac{s}{2} \quad ext{and} \quad \|oldsymbol{G}_{ au imes au}\| \leq rac{2 ext{K}_{ ext{G}}}{s} \, \|oldsymbol{G}\|_{\infty o 1} \, .$$

Since $2K_G < 4$ and $\boldsymbol{H}_{\tau \times \tau} = \boldsymbol{G}_{\tau \times \tau}$, we determine that

$$|\tau| \geq \frac{\mathrm{c}}{2} \cdot \mathrm{st.} \, \mathrm{rank}(\boldsymbol{A}) \quad \mathrm{and} \quad \|\boldsymbol{H}_{\tau \times \tau}\| \leq 0.5.$$

In view of Proposition 5.1, we conclude $\kappa(\mathbf{A}_{\tau}) \leq \sqrt{3}$.

Now, take another step back and notice that this here argument is nearly algorithmic. The random selection of σ can easily be implemented in practice, even though the proof does not specify the value of c. Given a Grothendieck factorization G = DTD, it is straightforward to identify the subset τ . The challenge, as before, is to produce the factorization.

5.3. Grothendieck factorization via convex optimization. As with the Pietsch factorization, the Grothendieck factorization can be identified from the solution to a convex program.

Theorem 5.5. Suppose G is Hermitian. The factorization G = DTD satisfies $||T|| \leq \alpha$ if and only if D satisfies

$$\lambda_{\max} \begin{bmatrix} -\alpha \mathbf{D}^2 & \mathbf{G} \\ \mathbf{G} & -\alpha \mathbf{D}^2 \end{bmatrix} \le 0. \tag{5.1}$$

In particular, if no D verifies this bound, then no factorization G = DTD admits $||T|| < \alpha$.

Proof. To check the forward implication, we essentially repeat the argument we used in Theorem 3.1 for the Pietsch case. This reasoning yields the pair of relations

$$G - \alpha D^2 \leq 0$$
 and $-G - \alpha D^2 \leq 0$.

Together, these two relations are equivalent with (5.1) because

$$\begin{bmatrix} -\alpha \mathbf{D}^2 & \mathbf{G} \\ \mathbf{G} & -\alpha \mathbf{D}^2 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} \mathbf{I} & \mathbf{I} \\ -\mathbf{I} & \mathbf{I} \end{bmatrix}^* \begin{bmatrix} \mathbf{G} - \alpha \mathbf{D}^2 \\ & -\mathbf{G} - \alpha \mathbf{D}^2 \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{I} \\ -\mathbf{I} & \mathbf{I} \end{bmatrix}.$$

To prove the reverse implication, we assume that (5.1) holds. First, we must check that $d_{jj} = 0$ implies that $g_j = 0$. To verify this claim, observe that

$$0 \ge \begin{bmatrix} \alpha \\ \boldsymbol{g}_j \end{bmatrix}^* \begin{bmatrix} 0 & \boldsymbol{g}_j^* \\ \boldsymbol{g}_j & -\alpha \boldsymbol{D}^2 \end{bmatrix} \begin{bmatrix} \alpha \\ \boldsymbol{g}_j \end{bmatrix} = \alpha \left(2 \|\boldsymbol{g}_j\|_2^2 - \boldsymbol{g}_j^* \boldsymbol{D}^2 \boldsymbol{g}_j \right) \ge \alpha \|\boldsymbol{g}_j\|_2^2$$

because trace(D^2) = 1. Therefore, we may construct a Grothendieck factorization G = DTD with $||T|| \le \alpha$ by setting $T = D^{\dagger}GD^{\dagger}$.

This discussion leads us to frame the eigenvalue minimization problem

min
$$\lambda_{\max}\begin{bmatrix} -\alpha \mathbf{F} & \mathbf{G} \\ \mathbf{G} & -\alpha \mathbf{F} \end{bmatrix}$$
 subject to trace $(\mathbf{F}) = 1$, \mathbf{F} diagonal, $\mathbf{F} \ge \mathbf{0}$ (5.2)

Owing to Theorem 5.5, there is a factorization G = DTD with $||T|| \le \alpha$ if and only if the value of (5.2) is nonpositive.

As in Section 3.2, we can easily construct Grothendieck factorizations from (imprecise) solutions to the problem (5.2). The proof of Bourgain–Tzafriri suggests that an appropriate value for the parameter $\alpha = s/4$. Furthermore, we do not need to solve (5.2) to optimality to obtain the required information. Indeed, it suffices to produce a feasible point with an objective value of O(1).

To solve (5.2) in practice, we again propose the Entropic Mirror Descent algorithm [BT03]. Appendix B describes the application to this problem. To provide a concrete bound on the computational cost, we remark that, when A_{τ} has dimension $m \times s$, forming $G = A_{\tau}^* A_{\tau} - \mathbf{I}$ costs at most $O(s^2 m)$, and Alizadeh's interior-point method [Ali95] requires $\widetilde{O}(s^{3.5})$ time.

Remark 5.6. For symmetric G, Theorem 5.3 shows that the norm $||G||_{\infty \to 1}$ is approximated within a factor K_G by the least α for which (5.2) has a nonpositive value. A natural reformulation of (5.2) can identify this value of α automatically (cf. Section 3.3). For nonsymmetric G, similar optimization problems arise. These ideas yield new approximation algorithms for the $(\infty, 1)$ norm.

5.4. An algorithm for Bourgain-Tzafriri. We are prepared to state our algorithm for producing the set τ described by the Bourgain-Tzafriri theorem. The procedure appears as Algorithm 2 on page 11. Note the striking similarity with Algorithm 1. The following result describes the performance of the algorithm. We omit the proof, which parallels that of Theorem 4.1.

Theorem 5.7. Suppose A is an $m \times n$ standardized matrix. With probability at least 3/4, Algorithm 2 produces a set $\tau = \tau_{\star}$ of column indices for which

$$|\tau| \ge c \cdot \text{st. rank}(\mathbf{A})$$
 and $\kappa(\mathbf{A}_{\tau}) \le \sqrt{3}$.

The computational cost is bounded by $\widetilde{O}(|\tau|^2 m + |\tau|^{3.5})$.

6. Future Directions

After the initial work [BT87], additional research has clarified the role of the stable rank. We highlight a positive result of Vershynin [Ver01, Cor. 7.1] and a negative result of Szarek [Sza90, Thm. 1.2] which together imply that the stable rank describes *precisely* how large a well-conditioned column submatrix can in general exist. See [Ver01, Sec. 5] for a more detailed discussion.

Theorem 6.1 (Vershynin 2001). Fix $\varepsilon > 0$. For each matrix \mathbf{A} , there is a set τ of column indices for which

$$|\tau| \ge (1 - \varepsilon) \cdot \text{st. rank}(\boldsymbol{A})$$
 and $\kappa(\boldsymbol{A}_{\tau}) \le C(\varepsilon)$.

Theorem 6.2 (Szarek). There is a sequence $\{A(n)\}$ of matrices of increasing dimension for which

$$|\tau| = \text{st. rank}(\mathbf{A}) \implies \kappa(\mathbf{A}_{\tau}) = \omega(1).$$

Vershynin's proof constructs the set τ in Theorem 6.1 with a complicated iteration that interleaves the Kashin–Tzafriri theorem and the Bourgain–Tzafriri theorem. We believe that the argument can be simplified substantially and developed into a column selection algorithm. This achievement might lead to a new method for performing rank-revealing factorizations, which could have a significant impact on the practice of numerical linear algebra.

Algorithm 1: Constructive version of Kashin–Tzafriri Theorem

```
KT(\boldsymbol{A})
Input: Standardized matrix \boldsymbol{A} with n columns
Output: A subset \tau_{\star} of \{1, 2, \ldots, n\}
Description: Produces \tau_{\star} such that |\tau_{\star}| \geq \text{st. rank}(A)/2 and ||A_{\tau}|| \leq 15 w.p. 4/5
     \tau_{\star} = \{1\}
1
     for s = 4, 8, 16, \dots, n
2
          for k = 1, 2, 3, \dots, 8 \log_2 s
               \tau = \text{Norm-Reduce}(\boldsymbol{A}, s)
4
                if ||A_{\tau}|| \leq 15 then \tau_{\star} = \tau and break
5
          if |\tau_{\star}| < s then exit
6
NORM-REDUCE(A, s)
Input: Standardized matrix A with n columns, a parameter s
Output: A subset \tau of \{1, 2, \ldots, n\}
     Draw a uniformly random set \sigma with cardinality s from \{1, 2, \ldots, n\}
     Solve (3.2) with \mathbf{B} = \mathbf{A}_{\sigma} and \alpha = 8 \mathrm{K}_{\mathrm{P}} \sqrt{s} to obtain a factorization \mathbf{B} = T \mathbf{D}
     Return \tau = \{j \in \sigma : d_{jj}^2 \le 2/s\}
```

Algorithm 2: Constructive version of Bourgain-Tzafriri Theorem

```
BT(\boldsymbol{A})
Input: Standardized matrix A with n columns
Output: A subset \tau_{\star} of \{1, 2, \ldots, n\}
Description: Produces \tau_{\star} such that |\tau_{\star}| \geq \text{st. rank}(\mathbf{A})/2 and \kappa(\mathbf{A}_{\tau}) \leq \sqrt{3} w.p. 3/4
     \tau_{\star} = \{1\}
     for s = 4, 8, 16, \dots, n
          for k = 1, 2, 3, \dots, 8 \log_2 s
3
               \tau = \text{Cond-Reduce}(\boldsymbol{A}, s)
4
               if \kappa(\mathbf{A}_{\tau}) \leq \sqrt{3} then \tau_{\star} = \tau and break
5
6
          if |\tau_{\star}| < s then exit
Cond-Reduce(\boldsymbol{A}, s)
Input: Standardized matrix A with n columns, a parameter s
Output: A subset \tau of \{1, 2, \ldots, n\}
     Draw a uniformly random set \sigma with cardinality s from \{1, 2, \dots, n\}
    Solve (5.2) with G = A_{\sigma}^* A_{\sigma} - I and \alpha = s/4 to obtain factorization G = DTD
    Return \tau = \{j \in \sigma : d_{ij}^2 \le 2/s\}
```

APPENDIX A. RANDOM REDUCTION OF NORMS

How does the norm of a matrix change when we pass to a random submatrix? This question has great importance in modern functional analysis, but it also has implications for the design of algorithms. This appendix describes some general results on how random selection reduces the $(\infty, 2)$ norm and the $(\infty, 1)$ norm. We also specialize these results to the structured matrices that appear in the proofs of Theorem 1.1 and Theorem 1.2.

A.1. Random Coordinate Models. We begin with two standard models for selecting random submatrices, and we describe how these models are related for an important class of matrix norms.

A matrix norm is *monotonic* if the norm of a matrix exceeds the norm of every (rectangular) submatrix. More precisely, the norm $\|\cdot\|$ is monotonic if

$$|||PAP'||| \leq ||A||$$

for each matrix A and each pair P, P' of diagonal (i.e., coordinate) projectors. The basic example of a monotonic matrix norm is the natural norm on operators from ℓ_p to ℓ_q with p,q in $[1,\infty]$, which is defined as

$$\|A\|_{p \to q} = \max\{\|Ax\|_q : \|x\|_p = 1\}.$$

Fix a number δ in [0,1], and denote by P_{δ} a random $n \times n$ diagonal matrix where exactly $s = \lfloor \delta n \rfloor$ entries equal one and the rest equal zero. This matrix can be viewed as a projector onto a random set of s coordinates. Therefore, we may treat AP_{δ} as a random s-column submatrix of A by ignoring the zeroed columns. Although this model is conceptually appealing, it can be difficult to analyze because of the dependencies among coordinates.

Let us introduce a simpler model for selecting random coordinates. We denote by \mathbf{R}_{δ} a random $n \times n$ diagonal matrix whose entries are independent 0–1 random variables with common mean δ . This matrix is a projector onto a random set of coordinates with average cardinality δn .

There is a basic result connecting these two models. The statement here follows directly from the argument in [Tro08, Lem. 14].

Proposition A.1 (Poissonization). Let $\|\cdot\|$ be a monotonic matrix norm. For each matrix A with n columns, it holds that

$$\mathbb{E} \| \boldsymbol{A} \boldsymbol{P}_{\delta} \| \leq 2 \, \mathbb{E} \| \boldsymbol{A} \boldsymbol{R}_{\delta} \|.$$

For each $n \times n$ matrix \mathbf{H} , it holds that

$$\mathbb{E} \| \boldsymbol{P}_{\delta} \boldsymbol{H} \boldsymbol{P}_{\delta} \| < 2 \mathbb{E} \| \boldsymbol{R}_{\delta} \boldsymbol{H} \boldsymbol{R}_{\delta} \|.$$

A.2. Reduction of the $(\infty, 2)$ norm. We begin with a general result on the $(\infty, 2)$ norm of a uniformly random set of columns drawn from a fixed matrix. The basic argument appears already in the work of Bourgain and Tzafriri [BT91, Thm. 1.1], but modern proofs are a little simpler. (See [Ver06, Lem. 2.3], for example.) The version here offers especially good constants.

Theorem A.2. Fix $\delta \in [0,1]$, and suppose **A** is a matrix with n columns. Then

$$\mathbb{E} \left\| \boldsymbol{A} \boldsymbol{R}_{\delta} \right\|_{\infty \to 2} \leq \sqrt{2\delta(1-\delta)} \left\| \boldsymbol{A} \right\|_{F} + \delta \left\| \boldsymbol{A} \right\|_{\infty \to 2}.$$

We postpone the argument to the next section so we may note a corollary that appears as a key step in the proof of the Kashin–Tzafriri theorem.

Corollary A.3. Suppose **A** is a standardized matrix with n columns. Choose $s \leq \lceil 2 \operatorname{st.rank}(\mathbf{A}) \rceil$, and write $\delta = s/n$. Then

$$\mathbb{E} \|\mathbf{A}\mathbf{P}_{\delta}\|_{\infty \to 2} \leq 7\sqrt{s}.$$

Proof. Owing to the standardization, $1 \le \text{st. rank}(\mathbf{A}) = n/\|\mathbf{A}\|^2$. It follows that

$$\delta \leq \frac{2\operatorname{st.rank}(\boldsymbol{A}) + 1}{n} \leq \frac{3\operatorname{st.rank}(\boldsymbol{A})}{n} = \frac{3}{\|\boldsymbol{A}\|^2}.$$

Apply the Poissonization result, Proposition A.1, to see that

$$\mathbb{E} \|AP_{\delta}\|_{\infty \to 2} \leq 2 \mathbb{E} \|AR_{\delta}\|_{\infty \to 2}$$
.

Theorem A.2 yields

$$\mathbb{E} \|\boldsymbol{A}\boldsymbol{P}_{\delta}\|_{\infty \to 2} \leq 2\sqrt{2\delta} \|\boldsymbol{A}\|_{F} + 2\delta \|\boldsymbol{A}\|_{\infty \to 2}.$$

Since \boldsymbol{A} has n unit-norm columns, it holds that $\|\boldsymbol{A}\|_{\mathrm{F}} = \sqrt{n}$. We also have the general bound $\|\boldsymbol{A}\|_{\infty \to 2} \leq \sqrt{n} \|\boldsymbol{A}\|$. Therefore,

$$\mathbb{E} \left\| \boldsymbol{A} \boldsymbol{P}_{\delta} \right\|_{\infty \to 2} \leq 2\sqrt{2\delta n} + 2\delta\sqrt{n} \left\| \boldsymbol{A} \right\| = 2\sqrt{s} \left[\sqrt{2} + \sqrt{\delta} \left\| \boldsymbol{A} \right\| \right].$$

Introduce the bound on δ and make a numerical estimate to complete the proof.

A.3. **Proof of Theorem A.2.** We must bound the quantity

$$E = \mathbb{E} \| A R_{\delta} \|_{\infty \to 2}$$
.

It turns out that it is easier to work with the (2,1) norm, which is dual to the $(\infty,2)$ norm, because there are some special methods that apply. Rewrite the expression as

$$E = \mathbb{E} \| \boldsymbol{R}_{\delta} \boldsymbol{A}^* \|_{2 \to 1} = \mathbb{E} \max_{\|\boldsymbol{x}\|_2 = 1} \sum_{j=1}^n \delta_j |\langle \boldsymbol{a}_j, \ \boldsymbol{x} \rangle|$$

where $\{\delta_j\}$ is a sequence of independent 0–1 random variables with common mean δ . In the sequel, we simplify notation by omitting the restriction on the vector \boldsymbol{x} and the limits from the sum.

The next step is to center and symmetrize the selectors. First, add and subtract the mean of each term from the sum and use the subadditivity of the maximum to obtain

$$\begin{split} E &\leq \mathbb{E} \max_{\boldsymbol{x}} \sum_{j} (\delta_{j} - \delta) \left| \left\langle \boldsymbol{a}_{j}, \ \boldsymbol{x} \right\rangle \right| + \max_{\boldsymbol{x}} \sum_{j} \delta \left| \left\langle \boldsymbol{a}_{j}, \ \boldsymbol{x} \right\rangle \right| \\ &= \mathbb{E} \max_{\boldsymbol{x}} \sum_{j} (\delta_{j} - \delta) \left| \left\langle \boldsymbol{a}_{j}, \ \boldsymbol{x} \right\rangle \right| + \delta \left\| \boldsymbol{A}^{*} \right\|_{2 \to 1} \\ &= \mathbb{E} \max_{\boldsymbol{x}} \sum_{j} (\delta_{j} - \delta) \left| \left\langle \boldsymbol{a}_{j}, \ \boldsymbol{x} \right\rangle \right| + \delta \left\| \boldsymbol{A} \right\|_{\infty \to 2}. \end{split}$$

We focus on the first term, which we abbreviate by the letter F. Let $\{\delta'_j\}$ be an independent copy of the sequence $\{\delta_j\}$. Jensen's inequality allows that

$$F = \mathbb{E} \max_{\boldsymbol{x}} \sum_{j} (\delta_{j} - \mathbb{E} \, \delta'_{j}) \, |\langle \boldsymbol{a}_{j}, \, \boldsymbol{x} \rangle|$$

$$\leq \mathbb{E} \max_{\boldsymbol{x}} \sum_{j} (\delta_{j} - \delta'_{j}) \, |\langle \boldsymbol{a}_{j}, \, \boldsymbol{x} \rangle|.$$

Observe that $\{\delta_j - \delta'_j\}$ is a sequence of independent, symmetric random variables. Thus, we may multiply each one by a random sign without changing the expectation [LT91, Lem. 6.3]. That is,

$$F \leq \mathbb{E} \max_{\boldsymbol{x}} \sum_{j} \varepsilon_{j} (\delta_{j} - \delta'_{j}) |\langle \boldsymbol{a}_{j}, \boldsymbol{x} \rangle|$$

where $\{\varepsilon_j\}$ is a sequence of independent Rademacher (i.e., uniform ± 1) random variables.

Now, we invoke a specific type of Rademacher comparison [LT91, Thm. 4.12 et seq.] to remove the absolute values from the inner product:

$$F \leq \mathbb{E} \max_{\boldsymbol{x}} \sum_{j} \varepsilon_{j} (\delta_{j} - \delta'_{j}) \langle \boldsymbol{a}_{j}, \boldsymbol{x} \rangle = \mathbb{E} \max_{\boldsymbol{x}} \left\langle \sum_{j} \varepsilon_{j} (\delta_{j} - \delta'_{j}) \boldsymbol{a}_{j}, \boldsymbol{x} \right\rangle.$$

Since x ranges over the ℓ_2 unit sphere, we reach

$$F \leq \mathbb{E} \left\| \sum_{j} \varepsilon_{j} (\delta_{j} - \delta'_{j}) \boldsymbol{a}_{j} \right\|_{2}.$$

The remaining expectations are elementary. First, apply Hölder's inequality to obtain

$$F \leq \left(\mathbb{E} \left\| \sum_{j} \varepsilon_{j} (\delta_{j} - \delta'_{j}) \boldsymbol{a}_{j} \right\|_{2}^{2} \right)^{1/2}.$$

Compute the expectation with respect to $\{\varepsilon_j\}$ and then with respect to $\{\delta_j\}$ and $\{\delta'_j\}$.

$$F \leq \left(\mathbb{E}\sum_{j} (\delta_{j} - \delta_{j}')^{2} \|\boldsymbol{a}_{j}\|_{2}^{2}\right)^{1/2} = \left(\sum_{j} 2\delta(1 - \delta) \|\boldsymbol{a}_{j}\|_{2}^{2}\right)^{1/2} = \sqrt{2\delta(1 - \delta)} \|\boldsymbol{A}\|_{F}.$$

Introduce this bound on F into the bound on E to conclude that

$$\mathbb{E} \|\boldsymbol{A}\boldsymbol{P}_{\delta}\|_{\infty \to 2} \leq \sqrt{2\delta(1-\delta)} \|\boldsymbol{A}\|_{F} + \delta \|\boldsymbol{A}\|_{\infty \to 2}.$$

This is the advertised estimate.

A.4. Reduction of the $(\infty, 1)$ norm. The impact of random selection on the $(\infty, 1)$ norm has already received some attention in the theoretical computer science literature because of a connection with graph cuts. The following result of Rudelson and Vershynin contains detailed information on the $(\infty, 1)$ norm of a random principal submatrix. The statement involves an auxiliary norm

$$\left\|oldsymbol{H}
ight\|_{ ext{col}} = \sum
olimits_j \left\|oldsymbol{H} \mathbf{e}_j
ight\|_2,$$

where $\{\mathbf{e}_j\}$ is the set of standard basis vectors. In words, we sum the Euclidean norms of the columns of the matrix.

Theorem A.4 (Rudelson-Vershynin). Fix $\delta \in [0,1]$, and suppose **H** is an $n \times n$ matrix. Then

$$\mathbb{E} \left\| \boldsymbol{R}_{\delta} \boldsymbol{H} \boldsymbol{R}_{\delta} \right\|_{\infty \to 1} \le C \left[\delta^{2} \left\| \boldsymbol{H} - \operatorname{diag}(\boldsymbol{H}) \right\|_{\infty \to 1} + \delta^{3/2} \left(\left\| \boldsymbol{H} \right\|_{\operatorname{col}} + \left\| \boldsymbol{H}^{*} \right\|_{\operatorname{col}} \right) + \delta \left\| \operatorname{diag}(\boldsymbol{H}) \right\|_{\infty \to 1} \right].$$

Theorem A.4 is established with the same methods as Theorem A.2, along with an additional decoupling argument [BT91, Prop. 1.9]. We rely on the following corollary in our proof of the Bourgain–Tzafriri theorem.

Corollary A.5. Suppose A is an n-column standardized matrix with hollow Gram matrix H = A*A - I. Choose $s \leq [c \cdot st.rank(A)]$, and write $\delta = s/n$. Then

$$\mathbb{E} \| \boldsymbol{P}_{\delta} \boldsymbol{H} \boldsymbol{P}_{\delta} \|_{\infty \to 1} \leq \frac{s}{9}.$$

Proof. Suppose A is a standardized matrix with n columns, and define its $n \times n$ hollow Gram matrix H. Observe that the $(\infty, 1)$ norm of H satisfies the bound

$$\|\boldsymbol{H}\|_{\infty \to 1} \le n \|\boldsymbol{H}\| \le n \max\{\|\boldsymbol{A}^*\boldsymbol{A}\| - 1, 1\} \le n \|\boldsymbol{A}\|^2.$$

Meanwhile, the $\|\cdot\|_{col}$ norm satisfies

$$\left\| \boldsymbol{H} \right\|_{\mathrm{col}} < \left\| \boldsymbol{A}^* \boldsymbol{A} \right\|_{\mathrm{col}} = \sum_{j} \left\| \boldsymbol{A}^* \boldsymbol{a}_j \right\|_2 \le n \left\| \boldsymbol{A} \right\|.$$

These facts play a central role in the calculation.

To continue, invoke the Poissonization result, Proposition A.1, which yields

$$\mathbb{E} \| P_{\delta} H P_{\delta} \|_{\infty \to 1} \leq 2 \| R_{\delta} H R_{\delta} \|_{\infty \to 1}.$$

Theorem A.4 provides that

$$\mathbb{E} \left\| \boldsymbol{P}_{\delta} \boldsymbol{H} \boldsymbol{P}_{\delta} \right\|_{\infty \to 1} \le C \left[\delta^{2} \left\| \boldsymbol{H} \right\|_{\infty \to 1} + \delta^{3/2} \left\| \boldsymbol{H} \right\|_{\text{col}} \right]$$

where we have applied the facts that H is Hermitian and has a zero diagonal. The two norm bounds result in additional simplifications:

$$\|\boldsymbol{P}_{\delta}\boldsymbol{H}\boldsymbol{P}_{\delta}\|_{\infty\to 1} \leq C \left[\delta^{2}n \|\boldsymbol{A}\|^{2} + \delta^{3/2}n \|\boldsymbol{A}\|\right] = Cs \left[\delta \|\boldsymbol{A}\|^{2} + \delta^{1/2} \|\boldsymbol{A}\|\right].$$

Since A has unit-norm columns, st. rank $(A) = n/\|A\|^2$. As a result, $\delta = s/n \le c/\|A\|^2$. By fixing a sufficiently small constant c, we can ensure that

$$\|\boldsymbol{P}_{\!\delta}\boldsymbol{H}\boldsymbol{P}_{\!\delta}\|_{\infty\to 1} \leq \frac{s}{9},$$

the advertised bound.

APPENDIX B. ENTROPIC MIRROR DESCENT

The algorithms for the Kashin–Tzafriri theorem and the Bourgain–Tzafriri theorem both require the solution to a convex minimization problem over the probability simplex. It is important to have a practical algorithm for approaching these optimizations. To that end, we briefly describe a simple, elegant method called Entropic Mirror Descent [BT03]. We then explain how to apply this technique to the specific objective functions that arise in our work.

B.1. Convex analysis. Let \mathbb{E} be a Euclidean space, i.e., a vector space equipped with a real-linear inner product. Let Ω be a convex subset of \mathbb{E} , and consider a convex function $J:\Omega\to\mathbb{R}$. The subdifferential $\partial J(f)$ contains each vector $\boldsymbol{\theta}\in\mathbb{E}^*$ that satisfies the inequalities

$$J(\boldsymbol{h}) - J(\boldsymbol{f}) \ge \langle \boldsymbol{\theta}, \ \boldsymbol{h} - \boldsymbol{f} \rangle$$
 for all $\boldsymbol{h} \in \Omega$.

The elements of the subdifferential are called *subgradients*. They describe the directions and rates of ascent of the function J at the point f. When J is differentiable at f, the gradient is the unique subgradient.

The Lipschitz constant of the function J with respect to a norm $\|\cdot\|$ is defined to be the least number L for which

$$|J(\boldsymbol{h}) - J(\boldsymbol{f})| \le L \|\boldsymbol{h} - \boldsymbol{f}\|$$
 for all $\boldsymbol{h}, \boldsymbol{f} \in \Omega$.

It can be shown [Roc70, Thm. 24.7] that

$$L = \sup\{ \|\boldsymbol{\theta}\|_* : \boldsymbol{\theta} \in \partial J(\boldsymbol{f}), \ \boldsymbol{f} \in \Omega \}.$$

where $\|\cdot\|_*$ is the dual norm.

B.2. Interior subgradient methods. Consider the (nonsmooth) convex program

min
$$J(\mathbf{f})$$
 subject to $\mathbf{f} \in \Omega$.

Subgradient information can be used to solve this problem, but caution is necessary because the negative subgradient is not necessarily a direction of descent. As a result, subgradient methods are typically *nonmonotone*, which means that the value of the objective function can (and often will) increase. It is also common for subgradient methods to produce iterates outside the constraint set. The classical remedy is to project each iterate back onto the constraint set. This idea succeeds, but it leads to zigzagging phenomena.

Interior subgradient methods [BT03] are designed to eliminate some of the problematic behavior that projected subgradient methods exhibit. To develop an interior subgradient method, we need a divergence measure that is tailored to the constraint set. At each iteration, we perform two steps:

(1) At the current iterate f, compute a subgradient $\theta \in \partial J(f)$ to linearize the objective function:

$$J(\boldsymbol{h}) \approx J(\boldsymbol{f}) + \langle \boldsymbol{\theta}, \ \boldsymbol{h} - \boldsymbol{f} \rangle$$

(2) Penalize the linearization with the divergence $D(\cdot; \mathbf{f})$ from the current iterate, scaled by a (large) parameter β^{-1} . Minimize this auxiliary function to produce a new iterate \mathbf{f}' :

$$f' \in \underset{h \in \Omega}{\operatorname{arg \ min}} \left\{ J(f) + \langle \boldsymbol{\theta}, \ \boldsymbol{h} - \boldsymbol{f} \rangle + \beta^{-1} D(\boldsymbol{h}; \boldsymbol{f}) \right\}.$$

Algorithm 3: Entropic Mirror Descent

```
EMD(J, s, T)
Input: Objective function J, dimension s, number T of iterations
Output: Approximate minimizer f of J

1 f^{(1)} = s^{-1}\mathbf{e} { Initialize with uniform density }
2 for t = 1 to T
3 Find \theta \in \partial J(f^{(t)}) { Compute subgradient }
4 \beta = \sqrt{\frac{2\log s}{T\|\theta\|_{\infty}^2}} { Compute step size }
5 h = f^{(t)} \cdot \exp(-\beta \theta) { Reweight current iterate }
6 f^{(t+1)} = h/\operatorname{trace}(h) { Rescale to obtain next iterate }
7 \operatorname{end} for
8 Return f \in \operatorname{arg\,min}_t J(f^{(t)})
```

The divergence penalty serves two purposes. First, it ensures that the next iterate is close to the previous iterate, which is essential because the linearization is only useful locally. Second, it simultaneously prevents the iterates from getting too close to the boundary of the constraint set. With a careful choice of the parameter β , we can guarantee progress toward the optimum set, at least on average.

B.3. Optimization on the probability simplex. The Entropic Mirror Descent (EMD) algorithm of Beck and Teboulle [BT03] is a specific instance of the interior subgradient method that is designed for minimizing convex functions over the probability simplex, the set defined by

$$\Delta_s = \{ \boldsymbol{f} \in \mathbb{R}^s : \operatorname{trace}(\boldsymbol{f}) = 1, \ \boldsymbol{f} \geq \boldsymbol{0} \}.$$

A natural divergence measure for this set is the relative entropy function:

$$D(\boldsymbol{h}; \boldsymbol{f}) = \sum_{j=1}^{s} h_j \log \left(\frac{h_j}{f_j}\right).$$

An amazing feature of the resulting interior subgradient method is that the optimization in the second step has a closed form:

$$h_j = \frac{f_j \exp(-\beta \theta_j)}{\sum_i f_j \exp(-\beta \theta_j)}$$
 for $j = 1, 2, \dots, s$.

Algorithm 3 describes the procedure that arises from these choices. Beck and Teboulle have established an elegant efficiency estimate [BT03, Thm. 4.2] for this method.

Theorem B.1 (Efficiency of EMD). Let $J: \Delta_s \to \mathbb{R}$ be a convex function whose Lipschitz constant with respect to the ℓ_1 norm is L. The approximate minimizer f generated by Algorithm 3 satisfies

$$J(\boldsymbol{f}) - J(\boldsymbol{f_{\star}}) \le \sqrt{\frac{2L^2 \log s}{T}}$$

where f_{\star} is a minimizer of J.

Algorithm 3 succeeds with a wide range of step sizes. In particular, when the total number T of iterations is unknown, we may compute the step size using the current iteration number t:

$$\beta = \sqrt{\frac{2\log s}{t \|\boldsymbol{\theta}\|_{\infty}^2}}.$$

This choice increases the right-hand side of the efficiency estimate by a logarithmic factor.

B.4. Pietsch factorization via EMD. Suppose B is a matrix with s columns. We can rephrase the Pietsch factorization problem (3.2) as an optimization over the probability simplex. Define the linear operator

$$\mathrm{diag}: \mathbb{R}^s \to \mathbb{R}^{s \times s}$$

that maps vectors to diagonal matrices in the obvious way. We can write the convex program as

min
$$\lambda_{\max}(\mathbf{B}^*\mathbf{B} - \alpha^2 \operatorname{diag}(\mathbf{f}))$$
 subject to $\mathbf{f} \in \Delta_s$. (B.1)

Abbreviate the objective function $J: \Delta_s \to \mathbb{R}$. We can evidently apply EMD to complete the optimization once we find a way to compute subgradients.

We use methods from the convex analysis of Hermitian matrices to determine the subdifferential of the objective function [Lew96]. Let A be an Hermitian matrix. Then

$$\partial \lambda_{\max}(\mathbf{A}) = \operatorname{conv}\{\mathbf{u}\mathbf{u}^* : \mathbf{A}\mathbf{u} = \lambda_{\max}(\mathbf{A})\mathbf{u}, \|\mathbf{u}\|_2 = 1\}.$$

In words, the subdifferential of the maximum eigenvalue function at \boldsymbol{A} is the convex hull of all rank-one projectors whose range lies in the top eigenspace of \boldsymbol{A} . According to [Roc70, Thm. 23.9], we have

$$\partial J(\mathbf{f}) = \{-\alpha^2 \operatorname{diag}^*(\mathbf{\Theta}) : \mathbf{\Theta} \in \partial \lambda_{\max}(\mathbf{B}^*\mathbf{B} - \alpha^2 \operatorname{diag}(\mathbf{f}))\}.$$

where the adjoint map diag* : $\mathbb{R}^{s \times s} \to \mathbb{R}^s$ extracts the diagonal of a matrix. In particular, we may construct a subgradient $\boldsymbol{\theta} \in \partial J(\boldsymbol{f})$ from a normalized maximal eigenvector \boldsymbol{u} of the matrix $\boldsymbol{B}^*\boldsymbol{B} - \alpha^2 \operatorname{diag}(\boldsymbol{f})$ using the formula

$$\boldsymbol{\theta} = -\alpha^2 \operatorname{diag}^*(\boldsymbol{u}\boldsymbol{u}^*) = -\alpha^2 |\boldsymbol{u}|^2$$

where $|\cdot|^2$ denotes the componentwise squared magnitude of a vector.

In summary, we can evaluate the objective function J(f) and simultaneously obtain a subgradient $\theta \in \partial J(f)$ from an eigenvector calculation plus some lower-order operations. Note that the standard methods for producing a single eigenvector, such as the Lanczos algorithm and its variants [GVL96, Ch. 9], require access to the matrix only through its action on vectors. It is therefore preferable in some settings—for example, when B is sparse—not to form the matrix B^*B .

Eigenvector computation is a primitive in every numerical linear package, so it is reasonable to assume that high-precision eigenvectors are available. In any case, slight variants of EMD will work with approximate subgradients, provided they are computed to sufficient precision. A simple analysis supporting this claim does not seem to be available in the optimization literature, but see [Kal07, Ch. 6] for related work.

We can bound the Lipschitz constant of J with respect to the ℓ_1 norm just by considering subgradients of the form $\theta = -\alpha^2 |u|^2$ because their convex hull yields the complete subdifferential. Since the eigenvector u is normalized,

$$\|\boldsymbol{\theta}\|_{\infty} = \alpha^2 \max_j |u_j|^2 \le \alpha^2,$$

we determine that the Lipschitz constant $L \leq \alpha^2$. According to Theorem B.1, the EMD algorithm ostensibly requires $\widetilde{O}(\alpha^4)$ iterations to deliver a solution to (B.1) with constant precision. In practice, far fewer iterations suffice.

Remark B.2. The application of EMD to (B.1) closely resembles the multiplicative weights method [Kal07, Ch. 6] for solving the MAXCUT problem (3.4). Indeed, the two algorithms are substantially identical, except for the specific choice of step sizes and the method for constructing the final solution from the sequence of iterates. The efficiency estimates are also similar, except that the multiplicative weights method uses the widths of the constraints in lieu of the Lipschitz constant. EMD appears to be more effective in practice because it exploits the geometry of the problem more completely.

B.5. Grothendieck factorization via EMD. Suppose G is an $s \times s$ Hermitian matrix. The Grothendieck factorization problem (5.2) can be expressed as solving

$$\min \ \lambda_{\max} egin{bmatrix} -lpha \operatorname{diag}(m{f}) & m{G} \ m{G} & -lpha \operatorname{diag}(m{f}) \end{bmatrix} \qquad ext{subject to} \qquad m{f} \in \Delta_s.$$

Abbreviate the objective function $J:\Delta_s\to\mathbb{R}$. Once again, EMD is an appropriate technique.

We may obtain subgradients using the same methods as before. Compute a normalized, maximal eigenvector of the matrix:

$$\begin{bmatrix} -\alpha \operatorname{diag}(\boldsymbol{f}) & \boldsymbol{G} \\ \boldsymbol{G} & -\alpha \operatorname{diag}(\boldsymbol{f}) \end{bmatrix} \begin{bmatrix} \boldsymbol{u} \\ \boldsymbol{v} \end{bmatrix} = J(\boldsymbol{f}) \begin{bmatrix} \boldsymbol{u} \\ \boldsymbol{v} \end{bmatrix} \quad \text{where} \quad \|\boldsymbol{u}\|_2^2 + \|\boldsymbol{v}\|_2^2 = 1.$$

Then a subgradient $\theta \in \partial J(f)$ can be obtained from the formula

$$\boldsymbol{\theta} = -\alpha \left(|\boldsymbol{u}|^2 + |\boldsymbol{v}|^2 \right).$$

The Lipschitz constant $L \leq \alpha$, so the number of iterations of EMD is apparently $\widetilde{O}(\alpha^2)$. Of course, the eigenvector calculations can be streamlined by exploiting the structure of the matrix.

Acknowledgments

The author thanks Ben Recht for helpful discussions about eigenvalue minimization.

References

- [Ali95] F. Alizadeh. Interior-point methods in semidefinite programming with applications to combinatorial optimization. SIAM J. Optimization, 5(1):13–51, Feb. 1995.
- [AN04] N. Alon and A. Naor. Approximating the cut norm via Grothendieck's inequality. In *Proc. 36th Ann. ACM Symposium on Theory of Computing (STOC)*, pages 72–80, Chicago, 2004.
- [BDM08] C. Boutsidis, P. Drineas, and M. Mahoney. On selecting exactly k columns from a matrix. Submitted for publication, 2008.
- [BT87] J. Bourgain and L. Tzafriri. Invertibility of "large" submatrices with applications to the geometry of Banach spaces and harmonic analysis. *Israel J. Math.*, 57(2):137–224, 1987.
- [BT91] J. Bourgain and L. Tzafriri. On a problem of Kadison and Singer. J. reine angew. Math., 420:1–43, 1991.
- [BT03] A. Beck and M. Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Res. Lett.*, 31:167–175, 2003.
- [GE96] M. Gu and S. Eisenstat. Efficient algorithms for computing a strong rank-revealing QR factorization. SIAM J. Sci. Comput., 17(4):848–869, Jul. 1996.
- [GVL96] G. H. Golub and C. F. Van Loan. Matrix Computations. Johns Hopkins Univ. Press, 3rd edition, 1996.
- [GW95] M. X. Goemans and D. P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. J. Assoc. Comput. Mach., 42:1115–1145, 1995.
- [HJ85] R. A. Horn and C. R. Johnson. Matrix Analysis. Cambridge Univ. Press, 1985.
- [Kal07] S. Kale. Efficient algorithms using the multiplicative weights update method. Ph.D. dissertation, Computer Science Dept., Princeton Univ., Princeton, Nov. 2007.
- [Lew96] A. S. Lewis. Convex analysis on the Hermitian matrices. SIAM J. Optimization, 6:164–177, 1996.
- [LO96] A. S. Lewis and M. L. Overton. Eigenvalue optimization. Acta Numerica, 5:149–190, 1996.
- [LT91] M. Ledoux and M. Talagrand. Probability in Banach Spaces: Isoperimetry and Processes. Springer, 1991.
- [Pis86] G. Pisier. Factorization of linear operators and geometry of Banach spaces. Number 60 in CBMS Regional Conference Series in Mathematics. AMS, Providence, 1986. Reprinted with corrections, 1987.
- [Roc70] R. T. Rockafellar. Convex Analysis. Princeton Univ. Press, 1970.
- [Roh00] J. Rohn. Computing the norm $||A||_{\infty,1}$ is NP-hard. Linear and Multilinear Algebra, 47:195–204, 2000.
- [RV07] M. Rudelson and R. Vershynin. Sampling from large matrices: An approach through geometric functional analysis. J. Amer. Comput. Soc., 54(4):Article 21, pp. 1–19, Jul. 2007.
- [Sza90] S. Szarek. Spaces with large distance from ℓ_{∞}^n and random matrices. Amer. J. Math., 112(6):899–942, Dec. 1990.
- [Tro08] J. A. Tropp. On the linear independence of spikes and sines. J. Fourier Anal. Appl., 2008. To appear.
- [Ver01] R. Vershynin. Johns decompositions: Selecting a large part. Israel J. Math., 122:253–277, 2001.
- [Ver06] R. Vershynin. High Dimensional Probability, volume 51 of IMS Lecture Notes—Monograph Series, chapter Random sets of isomorphism of linear operators on Hilbert space, pages 148–154. IMS, 2006.