

Subadditivity of matrix φ -entropy and concentration of random matrices

Richard Y. Chen*

Joel A. Tropp†

Abstract

This paper considers a class of entropy functionals defined for random matrices, and it demonstrates that these functionals satisfy a subadditivity property. Several matrix concentration inequalities are derived as an application of this result.

Keywords: Concentration inequalities; entropy method; moment inequalities; random matrix; subadditivity; tensorization.

AMS MSC 2010: Primary 60B20; 60E15, Secondary 60G09; 60F10.

Submitted to EJP on August 13, 2013, final version accepted on February 26, 2014.

Supersedes arXiv:1308.2952.

1 Introduction and Related Work

Entropy and related functions quantify the uncertainty inherent in a probability distribution. Measures of entropy have the property that the total entropy of a “product” is bounded by the sum of the entropies of the “factors.” This fundamental fact is called *subadditivity of entropy*, or sometimes *tensorization*, and it drives many applications of entropy. The survey [Lie75] contains a discussion of subadditivity in statistical mechanics, and the monograph [RS13] describes examples in information theory. In this work, we focus on the role of subadditivity of entropy in probability.

1.1 Subadditivity and Concentration

A *concentration inequality* states that a random variable is unlikely to exhibit a significant deviation from its mean value. The current intuition holds that a random variable concentrates whenever it depends smoothly on many independent random variables [Tal96]. Ledoux [Led97, Led01] and Bobkov & Ledoux [BL98] initiated a line of research that uses methods based on entropy to derive concentration inequalities. A few of the many authors who have contributed include Massart [Mas00a, Mas00b], Rio [Rio01], Bousquet [Bou02], and Boucheron et al. [BLM03, BLM09]. See the book [BLM13] for a comprehensive treatment of this theory and its bibliography.

*California Institute of Technology, USA. E-mail: ycchen@caltech.edu

†California Institute of Technology, USA. E-mail: jtropp@cms.caltech.edu

Let us summarize the ideas that lead from entropy to concentration. In this setting, we define the *entropy functional* for each nonnegative, real random variable Z :

$$H(Z) := \mathbb{E}(Z \log Z) - (\mathbb{E} Z) \log(\mathbb{E} Z). \tag{1.1}$$

Heuristically, $H(Z)$ quantifies our uncertainty about the precise value of Z . We typically consider the situation where $Z = e^{\theta Y}$ for a zero-mean random variable Y . In this case, we have the identity

$$\log \mathbb{E} e^{\theta Y} = \theta \int_0^\theta \frac{H(e^{\beta Y})}{\mathbb{E} e^{\beta Y}} \cdot \frac{d\beta}{\beta^2}. \tag{1.2}$$

Through Markov's inequality, bounds on the left-hand side imply that Y takes a large value with exponentially small probability. Therefore, we might hope to invoke inequalities for the entropy functional H to analyze the fluctuations of Y .

Indeed, the entropy functional exhibits a subadditivity property that allows us to implement this program. Suppose that Z is a function of mutually independent random variables X_1, \dots, X_n . We can define conditional entropy functionals

$$H_i(Z) := \mathbb{E}_i(Z \log Z) - (\mathbb{E}_i Z) \log(\mathbb{E}_i Z)$$

where \mathbb{E}_i denotes the expectation with respect to X_i , holding X_j fixed for $j \neq i$. The conditional entropy H_i reflects the uncertainty about Z that is attributable to our lack of knowledge about X_i . Subadditivity is the nontrivial result that

$$H(Z) \leq \sum_{i=1}^n \mathbb{E}[H_i(Z)]. \tag{1.3}$$

In other words, our uncertainty about Z does not exceed the total (average) uncertainty due to each X_i individually. Combining the identity (1.2) and the subadditivity property (1.3) with bounds for the conditional entropy, we can establish exponential concentration inequalities for functions of independent random variables.

The idea of considering alternative forms of entropy can be traced at least as far as the work of Rényi [Rén61], Bregman [Brè66], and Csiszár [Csi72]. In the early 2000s, researchers [LO00, Cha04, BBLM05, Cha06] recognized that generalized entropy functionals can exhibit subadditivity properties similar with those of the entropy functional (1.1). Let $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}$ be a convex function. The φ -entropy functional is defined for each nonnegative random variable Z by the formula

$$H_\varphi(Z) := \mathbb{E} \varphi(Z) - \varphi(\mathbb{E} Z).$$

The functional (1.1) derives from the choice $\varphi : t \mapsto t \log t$. Under stringent conditions on φ , it can be shown that the φ -entropy functional satisfies a subadditivity property analogous with (1.3). In particular, the function $\varphi : t \mapsto t^p$ yields a subadditive φ -entropy when $1 \leq p \leq 2$, a fact that leads to polynomial concentration inequalities [BBLM05].

1.2 Subadditivity of Matrix Entropies

The purpose of this paper is to explore the subadditivity properties of entropy functionals defined on matrix-valued random variables. Let $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}$ be a convex function. For a positive-semidefinite (psd) random matrix \mathbf{Z} , we can consider the *matrix φ -entropy functional*

$$H_\varphi(\mathbf{Z}) := \mathbb{E} \bar{\text{tr}} \varphi(\mathbf{Z}) - \bar{\text{tr}} \varphi(\mathbb{E} \mathbf{Z})$$

where φ refers to a standard matrix function and $\bar{\text{tr}}$ denotes the normalized trace. See Section 2.1 for definitions. It may be helpful to note some alternative presentations

of the matrix φ -entropy. First, the expression has the same structure as the scalar entropy (1.1) because

$$H_\varphi(\mathbf{Z}) = \mathbb{E} \Phi(\mathbf{Z}) - \Phi(\mathbb{E} \mathbf{Z}) \quad \text{where } \Phi := \bar{\text{tr}} \varphi \text{ is convex.}$$

Second, we can decompose the matrix entropy as

$$H_\varphi(\mathbf{Z}) = [\mathbb{E} \bar{\text{tr}} \varphi(\mathbf{Z}) - \varphi(\mathbb{E} \bar{\text{tr}} \mathbf{Z})] + [\varphi(\bar{\text{tr}} \mathbb{E} \mathbf{Z}) - \bar{\text{tr}} \varphi(\mathbb{E} \mathbf{Z})].$$

In other words, the matrix entropy quantifies the total loss in two different averaging operations on the matrix.

This work contains two main contributions:

1. We develop conditions on φ which ensure that the matrix φ -entropy is subadditive.
2. We verify these conditions for the functions $\varphi : t \mapsto t \log t$ and $\varphi : t \mapsto t^p$ where $p \in [1, 2]$.

The arguments parallel the analysis of scalar φ -entropies in Boucheron et al. [BBLM05], but the technical difficulties are more formidable in the matrix setting.

There are several areas that may benefit from this investigation.

Random matrix theory In the scalar setting, subadditivity of φ -entropy leads to powerful concentration inequalities. The subadditivity of matrix φ -entropy allows us to adapt these arguments to obtain some concentration inequalities for random matrices. See Section 1.3 for more information.

Convex analysis We derive subadditivity of the matrix φ -entropy functional H_φ from its convexity properties; see Lemma 4.1 et seq. These results may be useful in other contexts. For example, the convexity of scalar φ -entropy plays a role in machine learning [RW11, Sec. 2.5 et seq.].

Operator theory To prove that specific examples of matrix φ -entropy are subadditive, we rely on sophisticated methods from operator theory. In return, the results here may be relevant for problems in operator theory.

Quantum theory In quantum statistical mechanics and quantum information theory, entropies are defined for positive-definite matrices. Subadditivity of the quantum relative entropy function plays an important role in these fields, and this same result is closely connected with subadditivity of the matrix entropy H_φ where $\varphi : t \mapsto t \log t$. As such, subadditivity of other matrix φ -entropies may be relevant for quantum theory.

1.2.1 Related Work

After this paper was written, we learned about a contemporary paper [Han13] of Hansen that contains subadditivity results like the ones here. We also mention the subsequent paper [HZ14], which builds on our work. Detailed references appear below.

For the function $\varphi : t \mapsto t \log t$, we are aware of other precedents for our subadditivity results. In this special case, the subadditivity of H_φ follows from a classical result [Lin73] after a moderate amount of argument. In quantum statistical mechanics, subadditivity of entropy refers to a specific type of inequality for partitioned quantum systems [LR73]; see the lecture notes [Car10] for a recent presentation of these ideas. The paper [HOZ01] contains another type of subadditivity bound.

Finally, we mention the concept of *free entropy*, which is the appropriate generalization of entropy in noncommutative probability. As with other entropy measures, free entropy is subadditive. Arguments based on free entropy can be used to study extremely large matrices that are unitarily invariant. See the survey [Voi02] for further details and references.

1.3 Matrix Concentration Inequalities

A *matrix concentration inequality* provides a bound on the spectral-norm deviation of a random matrix from its mean [AW02, Oli09, Tro11, Tro12b, Min12, MJC⁺12, PMT13]; see the survey [Tro12c] for an annotated bibliography. These results have already found applications in a wide range of areas, including random graph theory [Oli09, CCT12], randomized linear algebra [DZ11, Tro12a, BG12], and least-squares approximation [CDL13]. There is also a separate line of work that leads to remarkable concentration inequalities for the spectral measure of a random Hermitian matrix; see for example [AGZ10, Sec. 4.4].

In spite of these successes, one frequently encounters random matrices that do not submit to existing techniques. Therefore, the study of matrix concentration remains an active area of investigation. As discussed in Section 1.1, subadditivity of scalar φ -entropy leads to a variety of concentration inequalities [BLM03, BBLM05]. It is natural to ask whether subadditivity of matrix φ -entropy leads to new concentration inequalities for random matrices.

We show that it is indeed possible to adapt scalar arguments based on subadditivity of entropy to the matrix setting, and we obtain some interesting matrix concentration inequalities. On the other hand, this method is not as satisfying as some other approaches to matrix concentration because the resulting bounds involve artificial assumptions. In fact, the matrix concentration inequalities in this work are dominated by the results we can obtain using arguments based on exchangeable pairs [PMT13]. This fact suggests that the subadditivity properties of matrix φ -entropy do not fully capture the behavior of a random matrix.

2 Main Results

In this section, we lay out detailed definitions and statements of our main results on subadditivity of matrix φ -entropy and its application to prove concentration inequalities for random matrices.

2.1 Notation and Background

Let us instate some notation. The set \mathbb{R}_+ contains the nonnegative real numbers, and \mathbb{R}_{++} consists of all positive real numbers. We write \mathbb{M}^d for the complex Banach space of $d \times d$ complex matrices, equipped with the usual ℓ_2 operator norm $\|\cdot\|$. The *normalized trace* is the function

$$\bar{\text{tr}} B := \frac{1}{d} \sum_{j=1}^d b_{jj} \quad \text{for } B \in \mathbb{M}^d.$$

The theory can be developed using the standard trace, but additional complications arise.

The set \mathbb{H}^d refers to the real-linear subspace of $d \times d$ Hermitian matrices in \mathbb{M}^d . For a matrix $A \in \mathbb{H}^d$, we write $\lambda_{\min}(A)$ and $\lambda_{\max}(A)$ for the algebraic minimum and maximum eigenvalues. For each interval $I \subset \mathbb{R}$, we define the set of Hermitian matrices whose eigenvalues fall in that interval:

$$\mathbb{H}^d(I) := \{A \in \mathbb{H}^d : [\lambda_{\min}(A), \lambda_{\max}(A)] \subset I\}.$$

We also introduce the set \mathbb{H}_+^d of $d \times d$ positive-semidefinite matrices and the set \mathbb{H}_{++}^d of $d \times d$ positive-definite matrices. Curly inequalities refer to the positive-semidefinite order. For example, $\mathbf{A} \preceq \mathbf{B}$ means that $\mathbf{B} - \mathbf{A}$ is positive semidefinite.

Next, let us explain how to extend scalar functions to matrices. Recall that each Hermitian matrix $\mathbf{A} \in \mathbb{H}^d$ has a *spectral resolution*

$$\mathbf{A} = \sum_{i=1}^d \lambda_i \mathbf{P}_i, \tag{2.1}$$

where $\lambda_1, \dots, \lambda_d$ are the eigenvalues of \mathbf{A} . The matrices $\mathbf{P}_1, \dots, \mathbf{P}_d$ are orthogonal projectors that satisfy the orthogonality relations

$$\mathbf{P}_i \mathbf{P}_j = \delta_{ij} \mathbf{P}_j \quad \text{and} \quad \sum_{i=1}^d \mathbf{P}_i = \mathbf{I},$$

where δ_{ij} is the Kronecker delta and \mathbf{I} is the identity matrix. One obtains a standard matrix function by applying a scalar function to the spectrum of a Hermitian matrix.

Definition 2.1 (Standard Matrix Function). *Let $f : I \mapsto \mathbb{R}$ be a function on an interval I of the real line. Suppose that $\mathbf{A} \in \mathbb{H}^d(I)$ has the spectral decomposition (2.1). Then*

$$f(\mathbf{A}) := \sum_{i=1}^d f(\lambda_i) \mathbf{P}_i.$$

We use lowercase Roman and Greek letters to refer to standard matrix functions. When we apply a familiar real-valued function to an Hermitian matrix, we are referring to the associated standard matrix function. Bold capital letters such as \mathbf{Y}, \mathbf{Z} denote general matrix functions that are not necessarily standard.

2.2 Subadditivity of Matrix Entropies

In this section, we provide an overview of the theory of matrix φ -entropies. At a high level, our approach has a strong parallel with the work of Boucheron et al. [BBLM05]. Nevertheless, there are interesting differences between the scalar and the matrix setting.

2.2.1 The Class of Matrix Entropies

First, we carve out a class of standard matrix functions that we can use to construct matrix entropies with the same subadditivity properties as their scalar counterparts.

Definition 2.2 (Φ_d Function Class). *Let d be a natural number. The class Φ_d contains each function $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}$ that is either affine or else satisfies the following three conditions.*

1. φ is convex and continuous at zero.
2. φ has two continuous derivatives on \mathbb{R}_{++} .
3. Define $\psi(t) = \varphi'(t)$ for $t \in \mathbb{R}_{++}$. The derivative $D\psi$ of the standard matrix function $\psi : \mathbb{H}_{++}^d \rightarrow \mathbb{H}^d$ is an invertible linear operator on \mathbb{H}_{++}^d , and the map $\mathbf{A} \mapsto [D\psi(\mathbf{A})]^{-1}$ is concave with respect to the semidefinite order on operators.

The technical definitions that support requirement (3) appear in Section 3. For now, we just remark that the scalar equivalent of (3) is the statement that $t \mapsto [\varphi''(t)]^{-1}$ is concave on \mathbb{R}_{++} .

The class Φ_1 coincides with the Φ function class considered in [BBLM05]. It can be shown that $\Phi_{d+1} \subseteq \Phi_d$ for each natural number d , so it is appropriate to introduce the class of matrix entropies:

$$\Phi_\infty := \bigcap_{d=1}^\infty \Phi_d$$

This class consists of scalar functions that satisfy the conditions of Definition 2.2 for an arbitrary choice of dimension d . Note that Φ_∞ is a convex cone: it contains all positive multiples and all finite sums of its elements.

In contrast to the scalar setting, it is quite hard to determine what functions are contained in Φ_∞ . The main technical achievement of this paper is to demonstrate that the standard entropy and certain power functions belong to the matrix entropy class.

Theorem 2.3 (Elements of the Matrix Entropy Class). *The following functions are members of the Φ_∞ class.*

1. The standard entropy $t \mapsto t \log t$.
2. The power function $t \mapsto t^p$ for each $p \in [1, 2]$.

The proof of Theorem 2.3 appears in Section 6. The statement about classical entropy can be obtained from standard results in matrix theory after some argument, but the result for power functions demands new effort. In fact, the claim about the classical entropy follows from the result for power functions because of the representation $t \log t = \lim_{p \downarrow 1} (t^p - t)/(p - 1)$.

See the independent work [Han13, Sec. 4] for closely related material. Very recently, Hansen and Zhang have developed an elegant characterization of the matrix entropy class [HZ14].

2.2.2 Matrix φ -Entropy

For each function in the matrix entropy class, we can introduce a generalized entropy functional that measures the amount of fluctuation in a random matrix.

Definition 2.4 (Matrix φ -Entropy). *Let $\varphi \in \Phi_\infty$. Consider a random matrix \mathbf{Z} taking values in \mathbb{H}_+^d , and assume that $\mathbb{E} \|\mathbf{Z}\| < \infty$ and $\mathbb{E} \|\varphi(\mathbf{Z})\| < \infty$. The matrix φ -entropy functional H_φ is*

$$H_\varphi(\mathbf{Z}) := \mathbb{E} \bar{\text{tr}} \varphi(\mathbf{Z}) - \bar{\text{tr}} \varphi(\mathbb{E} \mathbf{Z}). \tag{2.2}$$

Similarly, the conditional matrix φ -entropy functional is

$$H_\varphi(\mathbf{Z} \mid \mathcal{F}) := \mathbb{E} [\bar{\text{tr}} \varphi(\mathbf{Z}) \mid \mathcal{F}] - \bar{\text{tr}} \varphi(\mathbb{E}[\mathbf{Z} \mid \mathcal{F}]),$$

where \mathcal{F} is a subalgebra of the master sigma algebra.

For each convex function φ , the trace function $\bar{\text{tr}} \varphi : \mathbb{H}_+^d \rightarrow \mathbb{R}$ is also convex [Car10, Sec. 2.2]. Therefore, Jensen's inequality implies that the matrix φ -entropy is nonnegative:

$$H_\varphi(\mathbf{Z}) \geq 0.$$

For concreteness, here are some basic examples of matrix φ -entropy functionals.

$$\begin{aligned} H_\varphi(\mathbf{Z}) &= \bar{\text{tr}} [\mathbb{E}(\mathbf{Z} \log \mathbf{Z}) - (\mathbb{E} \mathbf{Z}) \log(\mathbb{E} \mathbf{Z})] && \text{when } \varphi(t) = t \log t. \\ H_\varphi(\mathbf{Z}) &= \bar{\text{tr}} [\mathbb{E}(\mathbf{Z}^p) - (\mathbb{E} \mathbf{Z})^p] && \text{when } \varphi(t) = t^p \text{ for } p \in [1, 2]. \\ H_\varphi(\mathbf{Z}) &= 0 && \text{when } \varphi \text{ is affine.} \end{aligned}$$

2.2.3 Subadditivity of Matrix φ -Entropy

The key fact about matrix φ -entropies is that they satisfy a subadditivity property. Let $\mathbf{x} := (X_1, \dots, X_n)$ denote a vector of independent random variables taking values in a Polish space, and write \mathbf{x}_{-i} for the random vector obtained by deleting the i th entry of \mathbf{x} .

$$\mathbf{x}_{-i} := (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n).$$

Consider a positive-semidefinite random matrix \mathbf{Z} that can be expressed as a measurable function of the random vector \mathbf{x} .

$$\mathbf{Z} := \mathbf{Z}(X_1, \dots, X_n) \in \mathbb{H}_+^d.$$

We instate the integrability conditions $\mathbb{E} \|\mathbf{Z}\| < \infty$ and $\mathbb{E} \|\varphi(\mathbf{Z})\| < \infty$.

Theorem 2.5 (Subadditivity of Matrix φ -Entropy). *Fix a function $\varphi \in \Phi_\infty$. Under the prevailing assumptions,*

$$H_\varphi(\mathbf{Z}) \leq \sum_{i=1}^n \mathbb{E} [H_\varphi(\mathbf{Z} | \mathbf{x}_{-i})]. \quad (2.3)$$

Typically, we apply Theorem 2.5 by way of a corollary. Let X'_1, \dots, X'_n denote independent copies of X_1, \dots, X_n , and form the random matrix

$$\mathbf{Z}'_i := \mathbf{Z}(X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n) \in \mathbb{H}_+^d.$$

Then \mathbf{Z}'_i and \mathbf{Z} are independent and identically distributed, conditional on the sigma algebra generated by \mathbf{x}_{-i} . In particular, these two random matrices are exchangeable counterparts.

Corollary 2.6 (Entropy Bounds via Exchangeability). *Fix a function $\varphi \in \Phi_\infty$, and write $\psi = \varphi'$. With the prevailing notation,*

$$H_\varphi(\mathbf{Z}) \leq \frac{1}{2} \sum_{i=1}^n \mathbb{E} \operatorname{tr} \left[(\mathbf{Z} - \mathbf{Z}'_i)(\psi(\mathbf{Z}) - \psi(\mathbf{Z}'_i)) \right].$$

Theorem 2.5 and Corollary 2.6 are matrix counterparts of the foundational results from Boucheron et al. [BBLM05, Sec. 3], which establish that scalar φ -entropies satisfy a similar subadditivity property. We devote Section 4 to the proof of these results.

2.3 Some Matrix Concentration Inequalities

Using Corollary 2.6, we can derive concentration inequalities for random matrices. In contrast to some previous approaches to matrix concentration, we need to place some significant restrictions on the type of random matrices we consider.

Definition 2.7 (Invariance under Signed Permutation). *A random matrix $\mathbf{Y} \in \mathbb{H}^d$ is invariant under signed permutation if we have the equality of distribution*

$$\mathbf{Y} \sim \mathbf{\Pi}^* \mathbf{Y} \mathbf{\Pi} \quad \text{for each signed permutation } \mathbf{\Pi}.$$

A signed permutation $\mathbf{\Pi} \in \mathbb{M}^d$ is a matrix with the properties that (i) each row and each column contains exactly one nonzero entry and (ii) the nonzero entries only take values $+1$ and -1 .

In particular, consider a random matrix that is invariant under orthogonal conjugation:

$$\mathbf{Y} \sim \mathbf{U}^* \mathbf{Y} \mathbf{U} \quad \text{for each orthogonal matrix } \mathbf{U}.$$

A matrix that satisfies this condition always verifies the requirement of Definition 2.7. Many classical ensembles, such as the GOE, satisfy this orthogonal invariance condition. Similar remarks apply to random matrices that are invariant under unitary conjugation.

2.3.1 A Bounded Difference Inequality

Let us present an exponential tail bound for a random matrix whose distribution is invariant under signed permutation.

Theorem 2.8 (Bounded Differences). *Let $\mathbf{x} := (X_1, \dots, X_n)$ be a vector of independent random variables, and let $\mathbf{x}' := (X'_1, \dots, X'_n)$ be an independent copy of \mathbf{x} . Consider random matrices*

$$\begin{aligned} \mathbf{Y} &:= \mathbf{Y}(X_1, \dots, X_i, \dots, X_n) \in \mathbb{H}^d \quad \text{and} \\ \mathbf{Y}'_i &:= \mathbf{Y}(X_1, \dots, X'_i, \dots, X_n) \in \mathbb{H}^d \quad \text{for } i = 1, \dots, n. \end{aligned}$$

Assume that \mathbf{Y} is invariant under signed permutation and that $\|\mathbf{Y}\|$ is bounded almost surely. Introduce the variance measure

$$V_{\mathbf{Y}} := \sup \left\| \mathbb{E} \left[\sum_{i=1}^n (\mathbf{Y} - \mathbf{Y}'_i)^2 \mid \mathbf{x} \right] \right\|, \quad (2.4)$$

where the supremum occurs over all possible values of \mathbf{x} . For each $t \geq 0$,

$$\begin{aligned} \mathbb{P} \{ \lambda_{\max}(\mathbf{Y} - \mathbb{E} \mathbf{Y}) \geq t \} &\leq d \cdot e^{-t^2/(2V_{\mathbf{Y}})}, \quad \text{and} \\ \mathbb{P} \{ \lambda_{\min}(\mathbf{Y} - \mathbb{E} \mathbf{Y}) \leq -t \} &\leq d \cdot e^{-t^2/(2V_{\mathbf{Y}})}. \end{aligned}$$

Theorem 2.8 follows from Corollary 2.6 with the choice $\varphi(t) = t \log t$. See Section 7 for the proof. This result can be viewed as a type of matrix bounded difference inequality. Closely related inequalities already appear in the literature; see [Tro12b, Cor. 7.5], [MJC⁺12, Cor. 11.1], and [PMT13, Cor. 4.1]. In fact, Theorem 2.8 is dominated by [PMT13, Cor. 4.1], which is not restricted to random matrices that are invariant under signed permutation.

2.3.2 Example: Sample covariance matrices

It may be helpful to sketch a short example that indicates the scope of Theorem 2.8. Consider a random vector of the form

$$\mathbf{w} := (\varepsilon_1 W_1, \varepsilon_2 W_2, \dots, \varepsilon_p W_p)^*$$

where (W_k) is an exchangeable family of random variables and (ε_k) is a sequence of independent Rademacher random variables. We also require that the random vector is bounded: $\|\mathbf{w}\|^2 \leq B$.

Let $\mathbf{w}_1, \dots, \mathbf{w}_n$ be iid copies of \mathbf{w} , and consider the sample covariance matrix

$$\mathbf{Y} := \frac{1}{n} \sum_{i=1}^n \mathbf{w}_i \mathbf{w}_i^*.$$

Our assumptions on \mathbf{w} ensure that \mathbf{Y} is invariant under signed permutation and that $\|\mathbf{Y}\|$ is bounded. Note that $\mathbb{E} \mathbf{Y} = c\mathbf{I}$ for a positive number c . It is also easy to check that the variance measure (2.4) satisfies $V_{\mathbf{Y}} \leq 2B^2/n$. An application of Theorem 2.8 delivers

$$\mathbb{P} \{ \|\mathbf{Y} - c\mathbf{I}\| \geq t \} \leq 2d \cdot e^{-nt^2/(4B^2)}.$$

The bound is informative when $c^2 > t^2 > 4B^2 \log(2d)/n$. In other words, the number n of samples should satisfy $n > 4B^2 \log(2d)/c^2$. Modulo constants, this estimate cannot be improved when \mathbf{w} has the uniform distribution on $\{\pm \mathbf{e}_1, \dots, \pm \mathbf{e}_p\}$, the set of signed standard basis vectors.

The main result of Rudelson's paper [Rud99] is a concentration bound for sample covariance matrices based on the noncommutative Khintchine inequality [LP86]. Rudelson allows any bounded random vector \mathbf{w} with a scalar covariance matrix, and he achieves the same result derived here.

2.3.3 Matrix Moment Bounds

We can also establish moment inequalities for a random matrix whose distribution is invariant under signed permutation.

Theorem 2.9 (Matrix Moment Bound). *Fix a number $q \in \{2, 3, 4, \dots\}$. Let $\mathbf{x} := (X_1, \dots, X_n)$ be a vector of independent random variables, and let $\mathbf{x}' := (X'_1, \dots, X'_n)$ be an independent copy of \mathbf{x} . Consider positive-semidefinite random matrices*

$$\begin{aligned} \mathbf{Y} &:= \mathbf{Y}(X_1, \dots, X_i, \dots, X_n) \in \mathbb{H}_+^d \quad \text{and} \\ \mathbf{Y}'_i &:= \mathbf{Y}(X_1, \dots, X'_i, \dots, X_n) \in \mathbb{H}_+^d \quad \text{for } i = 1, \dots, n. \end{aligned}$$

Assume that \mathbf{Y} is invariant under signed permutation and that $\mathbb{E}(\|\mathbf{Y}\|^q) < \infty$. Suppose that there is a constant $c \geq 0$ with the property

$$\mathbf{V}_{\mathbf{Y}} := \mathbb{E} \left[\sum_{i=1}^n (\mathbf{Y} - \mathbf{Y}'_i)^2 \mid \mathbf{x} \right] \preceq c \mathbf{Y}. \tag{2.5}$$

Then the random matrix \mathbf{Y} satisfies the moment inequality

$$[\mathbb{E} \bar{\text{tr}}(\mathbf{Y}^q)]^{1/q} \leq \mathbb{E} \bar{\text{tr}} \mathbf{Y} + \frac{q-1}{2} \cdot c.$$

Theorem 2.9 follows from Corollary 2.6 with the choice $\varphi(t) = t^{q/(q-1)}$. See Section 8 for the proof. This result can be regarded as a matrix extension of a moment inequality for real random variables [BBLM05, Cor. 1]. The paper [PMT13] contains similar moment inequalities for random matrices that need not satisfy the condition of signed permutation invariance. See also [JX03, JX08, JZ11, MJC⁺12].

2.4 Generalized Subadditivity of Matrix φ -Entropy

Theorem 2.5 is the shadow of a more sophisticated subadditivity property. We outline the simplest form of this more general result. See the lecture notes of Carlen [Car10] for more background on the topics in this section.

We work in the $*$ -algebra \mathbb{M}^d of $d \times d$ complex matrices, equipped with the conjugate transpose operation $*$ and the normalized trace inner product $\langle \mathbf{A}, \mathbf{B} \rangle := \bar{\text{tr}}(\mathbf{A}^* \mathbf{B})$. We say that a subspace $\mathfrak{A} \subset \mathbb{M}^d$ is a $*$ -subalgebra when \mathfrak{A} contains the identity matrix, \mathfrak{A} is closed under matrix multiplication, and \mathfrak{A} is closed under conjugate transposition. In other terms, $\mathbf{I} \in \mathfrak{A}$ and $\mathbf{AB} \in \mathfrak{A}$ and $\mathbf{A}^* \in \mathfrak{A}$ whenever $\mathbf{A}, \mathbf{B} \in \mathfrak{A}$.

In this setting, there is an elegant notion of conditional expectation. The orthogonal projector $\mathbb{E}_{\mathfrak{A}} : \mathbb{M}^d \rightarrow \mathfrak{A}$ onto the $*$ -subalgebra \mathfrak{A} is called the *conditional expectation* with respect to the $*$ -subalgebra. For $*$ -subalgebras \mathfrak{A} and \mathfrak{B} , we say that the conditional expectations $\mathbb{E}_{\mathfrak{A}}$ and $\mathbb{E}_{\mathfrak{B}}$ commute when

$$(\mathbb{E}_{\mathfrak{A}} \mathbb{E}_{\mathfrak{B}})(\mathbf{M}) = (\mathbb{E}_{\mathfrak{B}} \mathbb{E}_{\mathfrak{A}})(\mathbf{M}) \quad \text{for every } \mathbf{M} \in \mathbb{M}^d.$$

This construction generalizes the concept of independence in a probability space.

We can define the matrix φ -entropy conditional on a $*$ -subalgebra \mathfrak{A} :

$$H_{\varphi}(\mathbf{A} \mid \mathfrak{A}) := \bar{\text{tr}}[\varphi(\mathbf{A}) - \varphi(\mathbb{E}_{\mathfrak{A}} \mathbf{A})] \quad \text{for } \mathbf{A} \in \mathbb{H}_+^d.$$

Note that $\bar{\text{tr}}(\mathbb{E}_{\mathfrak{A}} \mathbf{A}) = \bar{\text{tr}} \mathbf{A}$ for each matrix \mathbf{A} in \mathbb{H}_+^d , so we do not need to include a conditional expectation in the leading term. Let $\mathfrak{A}_1, \dots, \mathfrak{A}_n$ be $*$ -subalgebras whose conditional expectations commute. Then we can extend the definition of the matrix φ -entropy to read

$$H_{\varphi}(\mathbf{A} \mid \mathfrak{A}_1, \dots, \mathfrak{A}_n) := \bar{\text{tr}}[\varphi(\mathbf{A}) - \varphi(\mathbb{E}_{\mathfrak{A}_1} \cdots \mathbb{E}_{\mathfrak{A}_n} \mathbf{A})] \quad \text{for } \mathbf{A} \in \mathbb{H}_+^d.$$

Because of commutativity, the order of the conditional expectations has no effect on the calculation. It turns out that matrix φ -entropy admits the following subadditivity property.

Theorem 2.10 (Subadditivity of Matrix φ -Entropy II). *Fix a function $\varphi \in \Phi_\infty$. Let $\mathfrak{A}_1, \dots, \mathfrak{A}_n$ be $*$ -subalgebras of \mathbb{M}^d whose conditional expectations commute. Then*

$$H_\varphi(\mathbf{A} | \mathfrak{A}_1, \dots, \mathfrak{A}_n) \leq \sum_{i=1}^n H_\varphi(\mathbf{A} | \mathfrak{A}_i) \quad \text{for } \mathbf{A} \in \mathbb{H}_+^d.$$

We omit the proof of this result. The argument involves considerations similar with Theorem 2.5, but it requires an extra dose of operator theory. The work in this paper already addresses the more challenging aspects of the proof. Note that the case $\varphi : t \mapsto t \log t$ is essentially a consequence of the classical results in [Lin73].

Theorem 2.10 can be seen as a formal extension of the subadditivity of matrix φ -entropy expressed in Theorem 2.5. To see why, let $\Omega := \Omega_1 \times \dots \times \Omega_n$ be a product probability space. The space $L_2(\Omega; \mathbb{M}^d)$ of random matrices is a $*$ -algebra with the normalized trace functional $\mathbb{E} \bar{\text{tr}}$. For each $i = 1, \dots, n$, we can form a $*$ -subalgebra \mathfrak{A}_i consisting of the random matrices that do not depend on the i th factor Ω_i of the product. The conditional expectation $\mathbb{E}_{\mathfrak{A}_i}$ simply integrates out the i th random variable. By independence, the family of conditional expectations $\mathbb{E}_{\mathfrak{A}_1}, \dots, \mathbb{E}_{\mathfrak{A}_n}$ commutes. Using this dictionary, compare the statement of Theorem 2.10 with Theorem 2.5.

3 Operators and Functions acting on Matrices

This work involves a substantial amount of operator theory. This section contains a short treatment of the basic facts. See [Bha97, Bha07] for a more complete introduction.

3.1 Linear Operators on Matrices

Let \mathbb{C}^d be the complex Hilbert space of dimension d , equipped with the standard inner product $\langle \mathbf{a}, \mathbf{b} \rangle := \mathbf{a}^* \mathbf{b}$. We usually identify \mathbb{M}^d with $\mathbb{B}(\mathbb{C}^d)$, the complex Banach space of linear operators acting on \mathbb{C}^d , equipped with the ℓ_2 operator norm $\|\cdot\|$.

We can also endow \mathbb{M}^d with the normalized trace inner product $\langle \mathbf{A}, \mathbf{B} \rangle := \bar{\text{tr}}(\mathbf{A}^* \mathbf{B})$ to form a Hilbert space. As a Hilbert space, \mathbb{M}^d is isometrically isomorphic with \mathbb{C}^{d^2} . Let $\mathbb{B}(\mathbb{M}^d)$ denote the complex Banach space of linear operators that map the Hilbert space \mathbb{M}^d into itself, equipped with the induced operator norm. The Banach space $\mathbb{B}(\mathbb{M}^d)$ is isometrically isomorphic with the Banach space \mathbb{M}^{d^2} .

As a consequence of this construction, every concept from matrix analysis has an immediate analog for linear operators on matrices. An operator $\mathsf{T} \in \mathbb{B}(\mathbb{M}^d)$ is *self-adjoint* when

$$\langle \mathbf{A}, \mathsf{T}(\mathbf{B}) \rangle = \langle \mathsf{T}(\mathbf{A}), \mathbf{B} \rangle \quad \text{for all } \mathbf{A}, \mathbf{B} \in \mathbb{B}(\mathbb{M}^d).$$

A self-adjoint operator $\mathsf{T} \in \mathbb{B}(\mathbb{M}^d)$ is *positive semidefinite* when

$$\langle \mathbf{A}, \mathsf{T}(\mathbf{A}) \rangle \geq 0 \quad \text{for all } \mathbf{A} \in \mathbb{M}^d.$$

For self-adjoint operators $\mathsf{S}, \mathsf{T} \in \mathbb{B}(\mathbb{M}^d)$, the notation $\mathsf{S} \preceq \mathsf{T}$ means that $\mathsf{T} - \mathsf{S}$ is positive semidefinite.

Each self-adjoint matrix operator $\mathsf{T} \in \mathbb{B}(\mathbb{M}^d)$ has a spectral resolution of the form

$$\mathsf{T} = \sum_{i=1}^{d^2} \lambda_i \mathsf{P}_i \tag{3.1}$$

where $\lambda_1, \dots, \lambda_{d^2}$ are the eigenvalues of T and the spectral projectors P_1, \dots, P_{d^2} are positive-semidefinite operators that satisfy

$$P_i P_j = \delta_{ij} P_j \quad \text{and} \quad \sum_{i=1}^{d^2} P_i = I,$$

where δ_{ij} is the Kronecker delta and I is the identity operator. As in the matrix case, a self-adjoint operator with nonnegative eigenvalues is the same thing as a positive-semidefinite operator.

We can extend a scalar function $f : I \rightarrow \mathbb{R}$ on an interval I of the real line to obtain a standard operator function. Indeed, if T has the spectral resolution (3.1) and the eigenvalues of T fall in the interval I , we define

$$f(T) := \sum_{i=1}^{d^2} f(\lambda_i) P_i.$$

This definition, of course, parallels the definition for matrices.

3.2 Monotonicity and Convexity

Let X and Y be sets of self-adjoint operators, such as $\mathbb{H}^d(I)$ or the set of self-adjoint operators in $\mathbb{B}(\mathbb{M}^d)$. We can introduce notions of monotonicity and convexity for a general function $\Psi : X \rightarrow Y$ using the semidefinite order on the spaces of operators.

Definition 3.1 (Monotone Operator-Valued Function). *The function $\Psi : X \rightarrow Y$ is monotone when*

$$S \preceq T \implies \Psi(S) \preceq \Psi(T) \quad \text{for all } S, T \in X.$$

Definition 3.2 (Convex Operator-Valued Function). *The function $\Psi : X \rightarrow Y$ is convex when X is a convex set and*

$$\Psi(\alpha S + \bar{\alpha} T) \preceq \alpha \cdot \Psi(S) + \bar{\alpha} \cdot \Psi(T) \quad \text{for all } \alpha \in [0, 1] \text{ and all } S, T \in X.$$

We have written $\bar{\alpha} := 1 - \alpha$. The function Ψ is concave when $-\Psi$ is convex.

The convexity of an operator-valued function Ψ is equivalent with a Jensen-type relation:

$$\Psi(\mathbb{E} X) \preceq \mathbb{E} \Psi(X) \tag{3.2}$$

whenever X is an integrable random operator taking values in X .

In particular, we can apply these definitions to standard matrix and operator functions. Let I be an interval of the real line. We say that the function $f : I \rightarrow \mathbb{R}$ is *operator monotone* when the lifted map $f : \mathbb{H}^d(I) \rightarrow \mathbb{H}^d$ is monotone for each natural number d . Likewise, the function $f : I \rightarrow \mathbb{R}$ is *operator convex* when the lifted map $f : \mathbb{H}^d(I) \rightarrow \mathbb{H}^d$ is convex for each natural number d .

Although scalar monotonicity and convexity are quite common, they are much rarer in the matrix setting [Bha97, Chap. 4]. For present purposes, we note that the power functions $t \mapsto t^p$ with $p \in [0, 1]$ are operator monotone and operator concave. The power functions $t \mapsto t^p$ with $p \in [1, 2]$ and the standard entropy $t \mapsto t \log t$ are all operator convex.

3.3 The Derivative of a Vector-Valued Function

The definition of the Φ_∞ function class involves a requirement that a certain standard matrix function is differentiable. For completeness, we include the background needed to interpret this condition.

Definition 3.3 (Derivative of a Vector-Valued Function). *Let X and Y be Banach spaces, and let U be an open subset of X . A function $F : U \rightarrow Y$ is differentiable at a point $A \in U$ if there exists a bounded linear operator $\mathsf{T} : X \rightarrow Y$ for which*

$$\lim_{\mathbf{B} \rightarrow 0} \frac{\|F(A + \mathbf{B}) - F(A) - \mathsf{T}(\mathbf{B})\|_Y}{\|\mathbf{B}\|_X} = 0.$$

When F is differentiable at A , the operator T is called the derivative of F at A , and we define $\mathsf{D}F(A) := \mathsf{T}$.

The derivative and the directional derivative have the following relationship:

$$\left. \frac{d}{ds} F(A + s\mathbf{B}) \right|_{s=0} = \mathsf{D}F(A)(\mathbf{B}). \tag{3.3}$$

In Section 6.2, we present an explicit formula for the derivative of a standard matrix function.

4 Subadditivity of Matrix φ -Entropy

In this section, we establish Theorem 2.5, which states that the matrix φ -entropy is subadditive for every function in the Φ_∞ class. This result depends on a variational representation for the matrix φ -entropy that appears in Section 4.1. We use the variational formula to derive a Jensen-type inequality in Section 4.2. The proof of Theorem 2.5 appears in Section 4.3.

4.1 Representation of Matrix φ -Entropy as a Supremum

The fundamental fact behind the subadditivity theorem is a representation of the matrix φ -entropy as a supremum of affine functions.

Lemma 4.1 (Supremum Representation for Entropy). *Fix a function $\varphi \in \Phi_\infty$, and introduce the scalar derivative $\psi = \varphi'$. Suppose that \mathbf{Z} is a random positive-semidefinite matrix for which $\|\mathbf{Z}\|$ and $\|\varphi(\mathbf{Z})\|$ are integrable. Then*

$$H_\varphi(\mathbf{Z}) = \sup_{\mathbf{T}} \mathbb{E} \bar{\text{tr}} [(\psi(\mathbf{T}) - \psi(\mathbb{E} \mathbf{T}))(\mathbf{Z} - \mathbf{T}) + \varphi(\mathbf{T}) - \varphi(\mathbb{E} \mathbf{T})]. \tag{4.1}$$

The range of the supremum contains each random positive-definite matrix \mathbf{T} for which $\|\mathbf{T}\|$ and $\|\varphi(\mathbf{T})\|$ are integrable. In particular, the matrix φ -entropy H_φ can be written in the dual form

$$H_\varphi(\mathbf{Z}) = \sup_{\mathbf{T}} \mathbb{E} \bar{\text{tr}} [\Upsilon_1(\mathbf{T}) \cdot \mathbf{Z} + \Upsilon_2(\mathbf{T})], \tag{4.2}$$

where $\Upsilon_i : \mathbb{H}_+^d \rightarrow \mathbb{H}^d$ for $i = 1, 2$.

This result implies that H_φ is a convex function on the space of random positive-semidefinite matrices. The dual representation of H_φ is well suited for establishing a form of Jensen's inequality, Lemma 4.3, which is the main ingredient in the proof of the subadditivity property, Theorem 2.5.

It may be valuable to see some particular instances of the dual representation of the matrix φ -entropy:

$$H_\varphi(\mathbf{Z}) = \sup_{\mathbf{T}} \mathbb{E} \bar{\text{tr}} [(\log \mathbf{T} - \log(\mathbb{E} \mathbf{T})) \cdot \mathbf{Z}] \quad \text{when } \varphi(t) = t \log t.$$

$$H_\varphi(\mathbf{Z}) = \sup_{\mathbf{T}} \mathbb{E} \bar{\text{tr}} [p(\mathbf{T}^{p-1} - (\mathbb{E} \mathbf{T})^{p-1}) \cdot \mathbf{Z} - (p-1)(\mathbf{T}^p - (\mathbb{E} \mathbf{T})^p)] \quad \text{when } \varphi(t) = t^p \text{ for } p \in [1, 2].$$

The first formula is the matrix version of a well-known variational principle for the classical entropy, cf. [BBLM05, p. 525]. In the matrix setting, this result can be derived from the joint convexity of quantum relative entropy [Lin73].

4.1.1 The Convexity Lemma

To establish the variational formula, we require a convexity result for a quadratic form connected with the function φ .

Lemma 4.2. *Fix a function $\varphi \in \Phi_\infty$, and let $\psi = \varphi'$. Suppose that \mathbf{Y} is a random matrix taking values in \mathbb{H}_+^d , and let \mathbf{K} be a random matrix taking values in \mathbb{M}^d . Assume that $\|\mathbf{Y}\|$ and $\|\mathbf{K}\|$ are integrable. Then*

$$\mathbb{E} \langle \mathbf{K}, \mathsf{D}\psi(\mathbf{Y})(\mathbf{K}) \rangle \geq \langle \mathbb{E} \mathbf{K}, \mathsf{D}\psi(\mathbb{E} \mathbf{Y})(\mathbb{E} \mathbf{K}) \rangle$$

Proof. The proof hinges on a basic convexity property of quadratic forms. Define a map that takes a matrix \mathbf{A} in \mathbb{H}^d and a positive-definite operator T on \mathbb{M}^d to a nonnegative number:

$$\mathcal{Q} : (\mathbf{A}, \mathsf{T}) \mapsto \langle \mathbf{A}, \mathsf{T}^{-1}(\mathbf{A}) \rangle.$$

We assert that the function \mathcal{Q} is convex. Indeed, the same result is well known when \mathbf{A} and T are replaced by a vector and a positive-definite matrix [Bha07, Exer. 1.5.1], and the extension is immediate from the isometric isomorphism between operators and matrices.

Recall that the Φ_∞ class requires $\mathbf{A} \mapsto [\mathsf{D}\psi(\mathbf{A})]^{-1}$ to be a concave map on \mathbb{H}_{++}^d . With these observations at hand, we can make the following calculation:

$$\begin{aligned} \mathbb{E} \langle \mathbf{K}, \mathsf{D}\psi(\mathbf{Y})(\mathbf{K}) \rangle &= \mathbb{E} \langle \mathbf{K}, ([\mathsf{D}\psi(\mathbf{Y})]^{-1})^{-1}(\mathbf{K}) \rangle \\ &\geq \langle \mathbb{E} \mathbf{K}, (\mathbb{E}[\mathsf{D}\psi(\mathbf{Y})]^{-1})^{-1}(\mathbb{E} \mathbf{K}) \rangle \\ &\geq \langle \mathbb{E} \mathbf{K}, ([\mathsf{D}\psi(\mathbb{E} \mathbf{Y})]^{-1})^{-1}(\mathbb{E} \mathbf{K}) \rangle \\ &= \langle \mathbb{E} \mathbf{K}, \mathsf{D}\psi(\mathbb{E} \mathbf{Y})(\mathbb{E} \mathbf{K}) \rangle. \end{aligned}$$

We obtain the second relation when we apply Jensen's inequality to the convex function \mathcal{Q} . The third relation depends on the semidefinite Jensen inequality (3.2) for the concave function $\mathbf{A} \mapsto [\mathsf{D}\psi(\mathbf{A})]^{-1}$, coupled with the fact [Bha97, Prop. V.1.6] that the operator inverse reverses the semidefinite order. \square

4.1.2 Proof of Lemma 4.1

The argument parallels the proof of [BBLM05, Lem. 1]. We begin with some reductions. The case where φ is an affine function is immediate, so we may require the derivative $\psi = \varphi'$ to be non-constant. By approximation, we may also assume that the random matrix \mathbf{Z} is strictly positive definite.

[Indeed, since φ is continuous on \mathbb{R}_+ , the Dominated Convergence Theorem implies that the matrix φ -entropy H_φ is continuous on the set containing each positive-semidefinite random matrix \mathbf{Y} where $\|\mathbf{Y}\|$ and $\|\varphi(\mathbf{Y})\|$ are integrable. Therefore, we can approximate a positive-semidefinite random matrix \mathbf{Z} by a sequence $\{\mathbf{Y}_n\}$ of positive-definite random matrices where $\mathbf{Y}_n \rightarrow \mathbf{Z}$ and be confident that $H_\varphi(\mathbf{Y}_n) \rightarrow H_\varphi(\mathbf{Z})$.]

When $\mathbf{T} = \mathbf{Z}$, the argument of the supremum in (4.1) equals $H_\varphi(\mathbf{Z})$. Therefore, our burden is to verify the inequality

$$H_\varphi(\mathbf{Z}) \geq \mathbb{E} \bar{\text{tr}} [(\psi(\mathbf{T}) - \psi(\mathbb{E} \mathbf{T}))(\mathbf{Z} - \mathbf{T}) + \mathbb{E} \varphi(\mathbf{T}) - \varphi(\mathbb{E} \mathbf{T})] \tag{4.3}$$

for each random positive-definite matrix \mathbf{T} that satisfies the same integrability requirements as \mathbf{Z} . For simplicity, we assume that the eigenvalues of both \mathbf{Z} and \mathbf{T} are bounded and bounded away from zero. See Appendix A for the extension to the general case.

We use an interpolation argument to establish (4.3). Define the family of random matrices

$$\mathbf{T}_s := (1 - s) \cdot \mathbf{Z} + s \cdot \mathbf{T} \quad \text{for } s \in [0, 1].$$

Introduce the real-valued function

$$F(s) := \mathbb{E} \bar{\text{tr}} [(\psi(\mathbf{T}_s) - \psi(\mathbb{E} \mathbf{T}_s)) \cdot (\mathbf{Z} - \mathbf{T}_s)] + H_\varphi(\mathbf{T}_s).$$

Observe that $F(0) = H_\varphi(\mathbf{Z})$, while $F(1)$ coincides with the right-hand side of (4.3). Therefore, to establish (4.3), it suffices to show that the function $F(s)$ is weakly decreasing on the interval $[0, 1]$.

We intend to prove that $F'(s) \leq 0$ for $s \in [0, 1]$. Since $\mathbf{Z} - \mathbf{T}_s = -s \cdot (\mathbf{T} - \mathbf{Z})$, we can rewrite the function F in the form

$$F(s) = -s \cdot \mathbb{E} \bar{\text{tr}} [(\psi(\mathbf{T}_s) - \psi(\mathbb{E} \mathbf{T}_s)) \cdot (\mathbf{T} - \mathbf{Z})] + \mathbb{E} \bar{\text{tr}} [\varphi(\mathbf{T}_s) - \varphi(\mathbb{E} \mathbf{T}_s)]. \quad (4.4)$$

We differentiate the function F to obtain

$$F'(s) = -s \cdot \mathbb{E} \bar{\text{tr}} [\text{D}\psi(\mathbf{T}_s)(\mathbf{T} - \mathbf{Z}) \cdot (\mathbf{T} - \mathbf{Z})] + s \cdot \bar{\text{tr}} [\text{D}\psi(\mathbb{E} \mathbf{T}_s)(\mathbb{E}(\mathbf{T} - \mathbf{Z})) \cdot (\mathbb{E}(\mathbf{T} - \mathbf{Z}))] \\ - \mathbb{E} \bar{\text{tr}} [(\psi(\mathbf{T}_s) - \psi(\mathbb{E} \mathbf{T}_s)) \cdot (\mathbf{T} - \mathbf{Z})] + \mathbb{E} \bar{\text{tr}} [(\psi(\mathbf{T}_s) - \psi(\mathbb{E} \mathbf{T}_s)) \cdot (\mathbf{T} - \mathbf{Z})]. \quad (4.5)$$

To handle the first term in (4.4), we applied the product rule, the rule (3.3) for directional derivatives, and the expression $d\mathbf{T}_s/ds = \mathbf{T} - \mathbf{Z}$. We used the identity $\text{D} \text{tr} \varphi(\mathbf{A}) = \psi(\mathbf{A})$ to differentiate the second term. We also relied on the Dominated Convergence Theorem to pass derivatives through expectations, which is justified because φ and ψ are continuously differentiable on \mathbb{H}_{++}^d and the eigenvalues of the random matrices are bounded and bounded away from zero. Now, the last two terms in (4.5) cancel, and we can rewrite the first two terms using the trace inner product:

$$F'(s) = s \cdot [\langle (\mathbb{E}(\mathbf{T} - \mathbf{Z})), \text{D}\psi(\mathbb{E} \mathbf{T}_s)(\mathbb{E}(\mathbf{T} - \mathbf{Z})) \rangle - \mathbb{E} \langle (\mathbf{T} - \mathbf{Z}), \text{D}\psi(\mathbf{T}_s)(\mathbf{T} - \mathbf{Z}) \rangle].$$

Invoke Lemma 4.2 to conclude that $F'(s) \leq 0$ for $s \in [0, 1]$.

4.2 A Conditional Jensen Inequality

The variational inequality in Lemma 4.1 leads directly to a Jensen inequality for the matrix φ -entropy.

Lemma 4.3 (Conditional Jensen Inequality). *Fix a function $\varphi \in \Phi_\infty$. Suppose that (X_1, X_2) is a pair of independent random variables taking values in a Polish space, and let $\mathbf{Z} = \mathbf{Z}(X_1, X_2)$ be a random positive-semidefinite matrix for which $\|\mathbf{Z}\|$ and $\|\varphi(\mathbf{Z})\|$ are integrable. Then*

$$H_\varphi(\mathbb{E}_1 \mathbf{Z}) \leq \mathbb{E}_1 H_\varphi(\mathbf{Z} | X_1),$$

where \mathbb{E}_1 is the expectation with respect to the first variable X_1 .

Proof. Let \mathbb{E}_2 denote the expectation with respect to the second variable X_2 . The result is a simple consequence of the dual representation (4.2) of the matrix φ -entropy:

$$H_\varphi(\mathbb{E}_1 \mathbf{Z}) = \sup_{\mathbf{T}} \mathbb{E}_2 \bar{\text{tr}} [\Upsilon_1(\mathbf{T}(X_2)) \cdot (\mathbb{E}_1 \mathbf{Z}) + \Upsilon_2(\mathbf{T}(X_2))]. \quad (4.6)$$

We have written $\mathbf{T}(X_2)$ to emphasize that this matrix depends only on the randomness in X_2 . To control (4.6), we apply Fubini's theorem to interchange the order of \mathbb{E}_1 and

\mathbb{E}_2 , and then we exploit the convexity of the supremum to draw out the expectation \mathbb{E}_1 .

$$\begin{aligned} H_\varphi(\mathbb{E}_1 \mathbf{Z}) &= \sup_{\mathbf{T}} \mathbb{E}_1 \mathbb{E}_2 \bar{\text{tr}} [\boldsymbol{\Upsilon}_1(\mathbf{T}(X_2)) \cdot \mathbf{Z} + \boldsymbol{\Upsilon}_2(\mathbf{T}(X_2))] \\ &\leq \mathbb{E}_1 \sup_{\mathbf{T}} \mathbb{E}_2 \bar{\text{tr}} [\boldsymbol{\Upsilon}_1(\mathbf{T}(X_2)) \cdot \mathbf{Z} + \boldsymbol{\Upsilon}_2(\mathbf{T}(X_2))] \\ &= \mathbb{E}_1 \sup_{\mathbf{T}} \mathbb{E} [\bar{\text{tr}}[\boldsymbol{\Upsilon}_1(\mathbf{T}(X_2)) \cdot \mathbf{Z} + \boldsymbol{\Upsilon}_2(\mathbf{T}(X_2))] | X_1] \\ &= \mathbb{E}_1 H_\varphi(\mathbf{Z} | X_1). \end{aligned}$$

The last relation is the duality formula (4.2), applied conditionally. \square

4.3 Proof of Theorem 2.5

We are now prepared to establish the main result on subadditivity of matrix φ -entropy. This theorem is a direct consequence of the conditional Jensen inequality, Lemma 4.3. In this argument, we write \mathbb{E}_i for the expectation with respect to the variable X_i . Using the notation from Section 2.2.3, we see that $\mathbb{E}_i = \mathbb{E}[\cdot | \mathbf{x}_{-i}]$.

First, separate the matrix φ -entropy into two parts by adding and subtracting terms:

$$\begin{aligned} H_\varphi(\mathbf{Z}) &= \mathbb{E} \bar{\text{tr}} [\varphi(\mathbf{Z}) - \varphi(\mathbb{E}_1 \mathbf{Z}) + \varphi(\mathbb{E}_1 \mathbf{Z}) - \varphi(\mathbb{E} \mathbf{Z})]. \\ &= \mathbb{E} [\mathbb{E}_1 \bar{\text{tr}} [\varphi(\mathbf{Z}) - \varphi(\mathbb{E}_1 \mathbf{Z})]] + \mathbb{E} \bar{\text{tr}} [\varphi(\mathbb{E}_1 \mathbf{Z}) - \varphi(\mathbb{E} \mathbf{Z})]. \end{aligned} \quad (4.7)$$

We can rewrite this expression as

$$\begin{aligned} H_\varphi(\mathbf{Z}) &= \mathbb{E} H_\varphi(\mathbf{Z} | \mathbf{x}_{-1}) + H_\varphi(\mathbb{E}_1 \mathbf{Z}) \\ &\leq \mathbb{E} H_\varphi(\mathbf{Z} | \mathbf{x}_{-1}) + \mathbb{E}_1 H_\varphi(\mathbf{Z} | X_1). \end{aligned} \quad (4.8)$$

The inequality follows from Lemma 4.3 because $\mathbf{Z} = \mathbf{Z}(X_1, \mathbf{x}_{-1})$ where X_1 and \mathbf{x}_{-1} are independent random variables.

The first term on the right-hand side of (4.8) coincides with the first summand on the right-hand side of the subadditivity inequality (2.3). We must argue that the remaining summands are contained in the second term on the right-hand side of (4.8). Repeating the argument in the previous paragraph, conditioning on X_1 , we obtain

$$H_\varphi(\mathbf{Z} | X_1) \leq \mathbb{E} [H_\varphi(\mathbf{Z} | \mathbf{x}_{-2}) | X_1] + \mathbb{E}_2 H_\varphi(\mathbf{Z} | X_1, X_2).$$

Substituting this expression into (4.8), we obtain

$$H_\varphi(\mathbf{Z}) \leq \sum_{i=1}^2 \mathbb{E} H_\varphi(\mathbf{Z} | \mathbf{x}_{-i}) + \mathbb{E}_1 \mathbb{E}_2 H_\varphi(\mathbf{Z} | X_1, X_2).$$

Continuing in this fashion, we arrive at the subadditivity inequality (2.3):

$$H_\varphi(\mathbf{Z}) \leq \sum_{i=1}^n \mathbb{E} H_\varphi(\mathbf{Z} | \mathbf{x}_{-i}).$$

This completes the proof of Theorem 2.5.

5 Entropy Bounds via Exchangeability

In this section, we derive Corollary 2.6, which uses exchangeable pairs to bound the conditional entropies that appear in Theorem 2.5. This result follows from another variational representation of the matrix φ -entropy.

5.1 Representation of the Matrix φ -Entropy as an Infimum

In this section, we present another formula for the matrix φ -entropy.

Lemma 5.1 (Infimum Representation for Entropy). *Fix a function $\varphi \in \Phi_\infty$, and let $\psi = \varphi'$. Assume that \mathbf{Z} is a random positive-semidefinite matrix where $\|\mathbf{Z}\|$ and $\|\varphi(\mathbf{Z})\|$ are integrable. Then*

$$H_\varphi(\mathbf{Z}) = \inf_{\mathbf{A} \in \mathbb{H}_+^d} \mathbb{E} \bar{\text{tr}} [\varphi(\mathbf{Z}) - \varphi(\mathbf{A}) - (\mathbf{Z} - \mathbf{A}) \cdot \psi(\mathbf{A})]. \quad (5.1)$$

Let \mathbf{Z}' be an independent copy of \mathbf{Z} . Then

$$H_\varphi(\mathbf{Z}) \leq \frac{1}{2} \cdot \mathbb{E} \bar{\text{tr}} [(\mathbf{Z} - \mathbf{Z}')(\psi(\mathbf{Z}) - \psi(\mathbf{Z}'))]. \quad (5.2)$$

We require a familiar trace inequality [Car10, Thm. 2.11]. This bound simply restates the fact that a convex function lies above its tangents.

Proposition 5.2 (Klein's Inequality). *Let $f : I \rightarrow \mathbb{R}$ be a differentiable convex function on an interval I of the real line. Then*

$$\bar{\text{tr}} [f(\mathbf{B}) - f(\mathbf{A}) - (\mathbf{B} - \mathbf{A}) \cdot f'(\mathbf{A})] \geq 0 \quad \text{for all } \mathbf{A}, \mathbf{B} \in \mathbb{H}^d(I).$$

With Klein's inequality at hand, the variational inequality follows quickly.

Proof of Lemma 5.1. Every function $\varphi \in \Phi_\infty$ is convex and differentiable, so Proposition 5.2 with $\mathbf{B} = \mathbb{E} \mathbf{Z}$ implies that

$$\bar{\text{tr}} [-\varphi(\mathbb{E} \mathbf{Z})] \leq \bar{\text{tr}} [-\varphi(\mathbf{A}) - (\mathbb{E} \mathbf{Z} - \mathbf{A}) \cdot \psi(\mathbf{A})]$$

for each fixed matrix $\mathbf{A} \in \mathbb{H}_+^d$. Substitute this bound into the definition (2.2) of the matrix φ -entropy, and draw the expectation out of the trace to reach

$$H_\varphi(\mathbf{Z}) \leq \mathbb{E} \bar{\text{tr}} [\varphi(\mathbf{Z}) - \varphi(\mathbf{A}) - (\mathbf{Z} - \mathbf{A}) \cdot \psi(\mathbf{A})]. \quad (5.3)$$

The inequality (5.3) becomes an equality when $\mathbf{A} = \mathbb{E} \mathbf{Z}$, which establishes the variational representation (5.1).

The symmetrized bound (5.2) follows from an exchangeability argument. Select $\mathbf{A} = \mathbf{Z}'$ in the expression (5.3), and apply the fact that $\mathbb{E} \varphi(\mathbf{Z}) = \mathbb{E} \varphi(\mathbf{Z}')$ to obtain

$$H_\varphi(\mathbf{Z}) \leq -\mathbb{E} \bar{\text{tr}} [(\mathbf{Z} - \mathbf{Z}') \cdot \psi(\mathbf{Z}')]. \quad (5.4)$$

Since \mathbf{Z} and \mathbf{Z}' are exchangeable, we can also bound the matrix φ -entropy as

$$H_\varphi(\mathbf{Z}) \leq -\mathbb{E} \bar{\text{tr}} [(\mathbf{Z}' - \mathbf{Z}) \cdot \psi(\mathbf{Z})]. \quad (5.5)$$

Take the average of the two bounds (5.4) and (5.5) to reach the desired inequality (5.2). \square

In the scalar case, stronger bounds are available. For a function $\varphi \in \Phi_1$,

$$\varphi(b) - \varphi(a) - (b - a)\varphi'(a) \leq (b - a)(\varphi'(b) - \varphi'(a)) \quad \text{for all } a, b \geq 0.$$

See [Cha06, Lem. 4.2] for details.

5.2 Proof of Corollary 2.6

Lemma 5.1 leads to a succinct proof of Corollary 2.6. We continue to use the notation from Section 2.2.3. Apply the inequality (5.2) conditionally to control the conditional matrix φ -entropy:

$$H_\varphi(\mathbf{Z} \mid \mathbf{x}_{-i}) \leq \frac{1}{2} \cdot \mathbb{E} \bar{\text{tr}} [(\mathbf{Z} - \mathbf{Z}'_i)(\psi(\mathbf{Z}) - \psi(\mathbf{Z}'_i)) \mid \mathbf{x}_{-i}] \quad (5.6)$$

because \mathbf{Z}'_i and \mathbf{Z} are conditionally iid, given \mathbf{x}_{-i} . Take the expectation on both sides of (5.6), and invoke the tower property of conditional expectation:

$$\mathbb{E} H_\varphi(\mathbf{Z} \mid \mathbf{x}_{-i}) \leq \frac{1}{2} \cdot \mathbb{E} \bar{\text{tr}} [(\mathbf{Z} - \mathbf{Z}'_i)(\psi(\mathbf{Z}) - \psi(\mathbf{Z}'_i))]. \quad (5.7)$$

To complete the proof, substitute (5.7) into the right-hand side of the bound (2.3) from the subadditivity result, Theorem 2.5.

6 Members of the Φ_∞ function class

In this section, we demonstrate that the classical entropy and certain power functions belong to the Φ_∞ function class. The main challenge is to verify that $\mathbf{A} \mapsto [\text{D}\psi(\mathbf{A})]^{-1}$ is a concave operator-valued map. We establish this result for the classical entropy in Section 6.4 and for the power function in Section 6.5. See the independent work [Han13, Sec. 4] for closely related results.

6.1 Tensor Product Operators

First, we explain the tensor product construction of an operator. The tensor product will allow us to represent the derivative of a standard matrix function compactly.

Definition 6.1 (Tensor Product). *Let $\mathbf{A}, \mathbf{B} \in \mathbb{H}^d$. The operator $\mathbf{A} \otimes \mathbf{B} \in \mathbb{B}(\mathbb{M}^d)$ is defined by the relation*

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{M}) = \mathbf{A}\mathbf{M}\mathbf{B} \quad \text{for each } \mathbf{M} \in \mathbb{M}^d. \quad (6.1)$$

The operator $\mathbf{A} \otimes \mathbf{B}$ is self-adjoint because we assume the factors are Hermitian matrices.

Suppose that $\mathbf{A}, \mathbf{B} \in \mathbb{H}^d$ are Hermitian matrices with spectral resolutions

$$\mathbf{A} = \sum_{i=1}^d \lambda_i \mathbf{P}_i \quad \text{and} \quad \mathbf{B} = \sum_{j=1}^d \mu_j \mathbf{Q}_j. \quad (6.2)$$

Then the tensor product $\mathbf{A} \otimes \mathbf{B}$ has the spectral resolution

$$\mathbf{A} \otimes \mathbf{B} = \sum_{i,j=1}^d \lambda_i \mu_j \mathbf{P}_i \otimes \mathbf{Q}_j.$$

In particular, the tensor product of two positive-definite matrices is a positive-definite operator.

6.2 The Derivative of a Standard Matrix Function

Next, we present some classical results on the derivative of a standard matrix function. See [Bha97, Sec. V.3] for further details.

Definition 6.2 (Divided Difference). *Let $f : I \rightarrow \mathbb{R}$ be a continuously differentiable function on an interval I of the real line. The first divided difference is the map $f^{[1]} : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by*

$$f^{[1]}(\lambda, \mu) := \begin{cases} f'(\lambda), & \lambda = \mu. \\ \frac{f(\lambda) - f(\mu)}{\lambda - \mu}, & \lambda \neq \mu, \end{cases}$$

We also require the Hermite representation of the divided difference:

$$f^{[1]}(\lambda, \mu) = \int_0^1 f'(\tau\lambda + \bar{\tau}\mu) d\tau, \tag{6.3}$$

where we have written $\bar{\tau} := 1 - \tau$.

The following result gives an explicit expression for the derivative of a standard matrix function in terms of a divided difference.

Proposition 6.3 (Daleckiĭ–Kreĭn Formula). *Let $f : I \rightarrow \mathbb{R}$ be a continuously differentiable function of an interval I of the real line. Suppose that $\mathbf{A} \in \mathbb{H}^d(I)$ is a diagonal matrix with $\mathbf{A} = \text{diag}(a_1, \dots, a_d)$. The derivative $Df(\mathbf{A}) \in \mathbb{B}(\mathbb{M}^d)$, and*

$$Df(\mathbf{A})(\mathbf{H}) = f^{[1]}(\mathbf{A}) \odot \mathbf{H} \quad \text{for } \mathbf{H} \in \mathbb{M}^d,$$

where \odot denotes the Schur (i.e., componentwise) product and $f^{[1]}(\mathbf{A})$ refers to the matrix of divided differences:

$$[f^{[1]}(\mathbf{A})]_{ij} = f^{[1]}(a_i, a_j) \quad \text{for } i, j = 1, \dots, d.$$

6.3 Operator Means

Our approach also relies on the concept of an operator mean. The following definition is due to Kubo & Ando [KA80].

Definition 6.4 (Operator Mean). *Let $f : \mathbb{R}_{++} \rightarrow \mathbb{R}_{++}$ be an operator concave function that satisfies $f(1) = 1$. Fix a natural number d . Let S and T be positive-definite operators in $\mathbb{B}(\mathbb{M}^d)$. We define the mean of the operators:*

$$M_f(S, T) := T^{1/2} \cdot f(T^{-1/2}ST^{-1/2}) \cdot T^{1/2} \in \mathbb{B}(\mathbb{M}^d).$$

When S and T commute, the formula simplifies to

$$M_f(S, T) = T \cdot f(ST^{-1}).$$

A few examples may be helpful. The function $f(s) = (1 + s)/2$ represents the usual arithmetic mean, the function $f(s) = s^{1/2}$ gives the geometric mean, and the function $f(s) = 2s/(1 + s)$ yields the harmonic mean. Operator means have a concavity property, which was established in the paper [KA80].

Proposition 6.5 (Operator Means are Concave). *Let $f : \mathbb{R}_{++} \rightarrow \mathbb{R}_{++}$ be an operator monotone function with $f(1) = 1$. Fix a natural number d . Suppose that S_1, S_2, T_1, T_2 are positive-definite operators in $\mathbb{B}(\mathbb{M}^d)$. Then*

$$\alpha \cdot M_f(S_1, T_1) + \bar{\alpha} \cdot M_f(S_2, T_2) \preceq M_f(\alpha S_1 + \bar{\alpha} S_2, \alpha T_1 + \bar{\alpha} T_2)$$

for $\alpha \in [0, 1]$ and $\bar{\alpha} = 1 - \alpha$.

6.4 Entropy

In this section, we demonstrate that the standard entropy function is a member of the Φ_∞ function class.

Theorem 6.6. *The function $\varphi : t \mapsto t \log t - t$ is a member of the Φ_∞ class.*

This result immediately implies Theorem 2.3(1), which states that $t \mapsto t \log t$ belongs to Φ_∞ . Indeed, the matrix entropy class contains all affine functions and all finite sums of its elements.

Theorem 6.6 follows easily from (deep) classical results because the variational representation of the standard entropy from Lemma 4.1 is equivalent with the joint convexity of quantum relative entropy [Lin73]. Instead of pursuing this idea, we present an argument that parallels the approach we use to study the power function. Some of the calculations below also appear in [Lie73, Proof of Cor. 2.1], albeit in compressed form.

Proof. Fix a positive integer d . We plan to show that the function $\varphi : t \mapsto t \log t - t$ is a member of the class Φ_d . Evidently, φ is continuous and convex on \mathbb{R}_+ , and it has two continuous derivatives on \mathbb{R}_{++} . It remains to verify the concavity condition for the second derivative.

Write $\psi(t) = \varphi'(t) = \log t$, and let $\mathbf{A} \in \mathbb{H}_{++}^d$. Without loss of generality, we may choose a basis where $\mathbf{A} = \text{diag}(a_1, \dots, a_d)$. The Daleckiĭ–Kreĭn formula, Proposition 6.3, tells us

$$D\psi(\mathbf{A})(\mathbf{H}) = \psi^{[1]}(\mathbf{A}) \odot \mathbf{H} = [\psi^{[1]}(a_i, a_j) \cdot h_{ij}]_{ij}.$$

As an operator, the derivative acts by Schur multiplication. This formula also makes it clear that the inverse of this operator acts by Schur multiplication:

$$[D\psi(\mathbf{A})]^{-1}(\mathbf{H}) = \left[\frac{1}{\psi^{[1]}(a_i, a_j)} \cdot h_{ij} \right]_{ij}.$$

Using the Hermite representation (6.3) of the first divided difference of $t \mapsto e^t$, we find

$$\frac{1}{\psi^{[1]}(\mu, \lambda)} = \frac{\lambda - \mu}{\log \lambda - \log \mu} = \int_0^1 e^{\tau \log \lambda + \bar{\tau} \log \mu} d\tau = \int_0^1 \lambda^\tau \mu^{\bar{\tau}} d\tau.$$

The latter calculation assumes that $\mu \neq \lambda$; it extends to the case $\mu = \lambda$ because both sides of the identity are continuous. As a consequence,

$$[D\psi(\mathbf{A})]^{-1}(\mathbf{H}) = \int_0^1 \left[a_i^{\bar{\tau}} h_{ij} a_j^{\tau} \right]_{ij} d\tau = \int_0^1 \mathbf{A}^\tau \mathbf{H} \mathbf{A}^{\bar{\tau}} d\tau = \int_0^1 (\mathbf{A}^\tau \otimes \mathbf{A}^{\bar{\tau}})(\mathbf{H}) d\tau.$$

We discover the expression

$$[D\psi(\mathbf{A})]^{-1} = \int_0^1 \mathbf{A}^\tau \otimes \mathbf{A}^{\bar{\tau}} d\tau. \tag{6.4}$$

This formula is correct for every positive-definite matrix.

For each $\tau \in [0, 1]$, consider the operator monotone function $f : t \mapsto t^\tau$ defined on \mathbb{R}_+ . Since $f(1) = 1$, we can construct the operator mean M_f associated with the function f . Note that $\mathbf{A} \otimes \mathbf{I}$ and $\mathbf{I} \otimes \mathbf{A}$ are commuting positive operators. Thus,

$$M_f(\mathbf{A} \otimes \mathbf{I}, \mathbf{I} \otimes \mathbf{A}) = (\mathbf{I} \otimes \mathbf{A}) \cdot f((\mathbf{A} \otimes \mathbf{I})(\mathbf{I} \otimes \mathbf{A})^{-1}) = \mathbf{A}^\tau \otimes \mathbf{A}^{\bar{\tau}}.$$

The map $\mathbf{A} \mapsto (\mathbf{A} \otimes \mathbf{I}, \mathbf{I} \otimes \mathbf{A})$ is linear, so Proposition 6.5 guarantees that $\mathbf{A} \mapsto \mathbf{A}^\tau \otimes \mathbf{A}^{\bar{\tau}}$ is concave for each $\tau \in [0, 1]$. This result is usually called the Lieb Concavity Theorem [Bha97, Thm. IX.6.1]. Combine this fact with the integral representation (6.4) to reach the conclusion that $\mathbf{A} \mapsto [D\psi(\mathbf{A})]^{-1}$ is a concave map on the cone \mathbb{H}_{++}^d of positive-definite matrices. \square

6.5 Power Functions

In this section, we prove that certain power functions belong to the Φ_∞ function class.

Theorem 6.7. *For each $p \in [0, 1]$, the function $\varphi : t \mapsto t^{p+1}/(p+1)$ is a member of the Φ_∞ class.*

This result immediately implies Theorem 2.3(2), which states that $t \mapsto t^{p+1}$ belongs to the class Φ_∞ . Indeed, the matrix entropy class contains all positive multiples of its elements.

The proof of Theorem 6.7 follows the same path as Theorem 6.6, but it is somewhat more involved. First, we derive an expression for the function $\mathbf{A} \mapsto [\mathbf{D}\psi(\mathbf{A})]^{-1}$ where $\psi = \varphi'$.

Lemma 6.8. *Fix $p \in (0, 1]$, and let $\psi(t) = t^p$ for $t \geq 0$. For each matrix $\mathbf{A} \in \mathbb{H}_+^d$,*

$$[\mathbf{D}\psi(\mathbf{A})]^{-1} = \frac{1}{p} \int_0^1 (\tau \cdot \mathbf{A}^p \otimes \mathbf{I} + \bar{\tau} \cdot \mathbf{I} \otimes \mathbf{A}^p)^{(1-p)/p} d\tau, \quad (6.5)$$

where $\bar{\tau} := 1 - \tau$.

Proof. As before, we may assume without loss of generality that the matrix $\mathbf{A} = \text{diag}(a_1, \dots, a_d)$. Using the Daleckiĭ–Kreĭn formula, Proposition 6.3, we see that

$$[\mathbf{D}\psi(\mathbf{A})]^{-1}(\mathbf{H}) = \left[\frac{1}{\psi^{[1]}(a_i, a_j)} \cdot h_{ij} \right].$$

The Hermite representation (6.3) of the first divided difference of $t \mapsto t^{1/p}$ gives

$$\frac{1}{\psi^{[1]}(\mu, \lambda)} = \frac{\mu - \lambda}{\mu^p - \lambda^p} = \frac{1}{p} \int_0^1 (\tau \cdot \lambda^p + \bar{\tau} \cdot \mu^p)^{(1-p)/p} d\tau =: g(\lambda, \mu).$$

We use continuity to verify that the latter calculation remains valid when $\mu = \lambda$. Using this function g , we can identify a compact representation of the operator:

$$[\mathbf{D}\psi(\mathbf{A})]^{-1}(\mathbf{H}) = \sum_{ij} g(a_i, a_j) h_{ij} \mathbf{E}_{ij} = \left[\sum_{ij} g(a_i, a_j) (\mathbf{E}_{ii} \otimes \mathbf{E}_{jj}) \right](\mathbf{H}),$$

where we write \mathbf{E}_{ij} for the matrix with a one in the (i, j) position and zeros elsewhere. It remains to verify that the bracket coincides with the expression (6.5). Indeed,

$$\begin{aligned} \sum_{ij} g(a_i, a_j) (\mathbf{E}_{ii} \otimes \mathbf{E}_{jj}) &= \frac{1}{p} \int_0^1 \sum_{ij} (\tau \cdot a_i^p + \bar{\tau} \cdot a_j^p)^{(1-p)/p} (\mathbf{E}_{ii} \otimes \mathbf{E}_{jj}) d\tau \\ &= \frac{1}{p} \int_0^1 \left[\sum_{ij} (\tau \cdot a_i^p + \bar{\tau} \cdot a_j^p) (\mathbf{E}_{ii} \otimes \mathbf{E}_{jj}) \right]^{(1-p)/p} d\tau \\ &= \frac{1}{p} \int_0^1 (\tau \cdot \mathbf{A}^p \otimes \mathbf{I} + \bar{\tau} \cdot \mathbf{I} \otimes \mathbf{A}^p)^{(1-p)/p} d\tau. \end{aligned}$$

The second relation follows from the definition of the standard operator function associated with $t \mapsto t^{(1-p)/p}$. To confirm that the third line equals the second, expand the matrices $\mathbf{A} = \sum_i a_i \mathbf{E}_{ii}$ and $\mathbf{I} = \sum_j \mathbf{E}_{jj}$ and invoke the bilinearity of the tensor product. \square

Proof of Theorem 6.7. We are now prepared to prove that certain power functions belong to the Φ_∞ function class. Fix an exponent $p \in [0, 1]$, and let d be a fixed positive

integer. We intend to show that the function $\varphi(t) = t^{p+1}/(p+1)$ belongs to the Φ_d class. When $p = 0$, the function φ is affine, so we may assume that $p > 0$. It is clear that φ is continuous and convex on \mathbb{R}_+ , and φ has two continuous derivatives on \mathbb{R}_{++} . It remains to verify that the second derivative has the required concavity property.

Let $\psi(t) = \varphi'(t) = t^p$ for $t \geq 0$, and consider a matrix $\mathbf{A} \in \mathbb{H}_{++}^d$. Lemma 6.8 demonstrates that

$$[\mathbf{D}\psi(\mathbf{A})]^{-1} = \frac{1}{p} \int_0^1 (\tau \cdot \mathbf{A}^p \otimes \mathbf{I} + \bar{\tau} \cdot \mathbf{I} \otimes \mathbf{A}^p)^{(1/p)(1-p)} d\tau, \tag{6.6}$$

where we maintain the usage $\bar{\tau} := 1 - \tau$. For each $\tau \in [0, 1]$, the scalar function $a \mapsto \tau a + \bar{\tau}$ is operator monotone because it is affine and increasing. On account of the result [And79, Cor. 4.3], the function

$$f : a \mapsto (\tau \cdot a^p + \bar{\tau})^{1/p}$$

is also operator monotone. A short calculation shows that $f(1) = 1$. Therefore, we can use f to construct an operator mean M_f . Since $\mathbf{A} \otimes \mathbf{I}$ and $\mathbf{I} \otimes \mathbf{A}$ are commuting positive operators, we have

$$M_f(\mathbf{A} \otimes \mathbf{I}, \mathbf{I} \otimes \mathbf{A}) = (\mathbf{I} \otimes \mathbf{A}) \cdot f((\mathbf{A} \otimes \mathbf{I})(\mathbf{I} \otimes \mathbf{A})^{-1}) = (\tau \cdot \mathbf{A}^p \otimes \mathbf{I} + \bar{\tau} \cdot \mathbf{I} \otimes \mathbf{A}^p)^{1/p}.$$

The map $\mathbf{A} \mapsto (\mathbf{A} \otimes \mathbf{I}, \mathbf{I} \otimes \mathbf{A})$ is linear, so Proposition 6.5 ensures that

$$\mathbf{A} \mapsto (\tau \cdot \mathbf{A}^p \otimes \mathbf{I} + \bar{\tau} \cdot \mathbf{I} \otimes \mathbf{A}^p)^{1/p} \tag{6.7}$$

is a concave map.

We are now prepared to check that (6.6) defines a concave operator. Let \mathbf{S}, \mathbf{T} be arbitrary positive-definite matrices, and choose $\alpha \in [0, 1]$. Suppose that \mathbf{Z} is the random matrix that takes value \mathbf{S} with probability α and value \mathbf{T} with probability $1 - \alpha$. For each $\tau \in [0, 1]$, we compute

$$\begin{aligned} \mathbb{E} [(\tau \cdot \mathbf{Z}^p \otimes \mathbf{I} + \bar{\tau} \cdot \mathbf{I} \otimes \mathbf{Z}^p)^{1/p}]^{1-p} &\preceq [\mathbb{E} (\tau \cdot \mathbf{Z}^p \otimes \mathbf{I} + \bar{\tau} \cdot \mathbf{I} \otimes \mathbf{Z}^p)^{1/p}]^{1-p} \\ &\preceq [(\tau \cdot (\mathbb{E} \mathbf{Z})^p \otimes \mathbf{I} + \bar{\tau} \cdot \mathbf{I} \otimes (\mathbb{E} \mathbf{Z})^p)^{1/p}]^{1-p}. \end{aligned}$$

The first relation holds because $t \mapsto t^{1-p}$ is operator concave [Bha97, Thm. V.1.9 and Thm. V.2.5]. To obtain the second relation, we apply the concavity property of the map (6.7), followed by the fact that $t \mapsto t^{1-p}$ is operator monotone [Bha97, Thm. V.1.9]. This calculation establishes the claim that

$$\mathbf{A} \mapsto (\tau \cdot \mathbf{A}^p \otimes \mathbf{I} + \bar{\tau} \cdot \mathbf{I} \otimes \mathbf{A}^p)^{(1-p)/p}$$

is concave on \mathbb{H}_{++}^d for each $\tau \in [0, 1]$. In view of the integral representation (6.6), we may conclude that $\mathbf{A} \mapsto [\mathbf{D}\psi(\mathbf{A})]^{-1}$ is concave on the cone \mathbb{H}_{++}^d of positive-definite matrices. \square

7 A Bounded Difference Inequality for Random Matrices

In this section, we prove Theorem 2.8, a bounded difference inequality for a random matrix whose distribution is invariant under signed permutation. We begin with some preliminaries that support the proof, and we establish the main result in Section 7.2.

7.1 Preliminaries

First, we describe how to compute the expectation of a function of a random matrix whose distribution is invariant under signed permutation. See Definition 2.7 for a reminder of what this requirement means.

Lemma 7.1. *Let $f : I \rightarrow \mathbb{R}$ be a function on an interval I of the real line. Assume that $\mathbf{X} \in \mathbb{H}^d(I)$ is a random matrix whose distribution is invariant under signed permutation. Then*

$$\mathbb{E} f(\mathbf{X}) = \bar{\text{tr}}[\mathbb{E} f(\mathbf{X})] \cdot \mathbf{I}.$$

Proof. Let $\mathbf{\Pi} \in \mathbb{H}^d$ be an arbitrary signed permutation matrix. Observe that

$$\mathbb{E} f(\mathbf{X}) = \mathbb{E} f(\mathbf{\Pi}^* \mathbf{X} \mathbf{\Pi}) = \mathbf{\Pi}^* [\mathbb{E} f(\mathbf{X})] \mathbf{\Pi}. \quad (7.1)$$

The first relation holds because the distribution of \mathbf{X} is invariant under conjugation by $\mathbf{\Pi}$. The second relation follows from the definition of a standard matrix function and the fact that $\mathbf{\Pi}$ is unitary. We may average (7.1) over $\mathbf{\Pi}$ drawn from the uniform distribution on the set of signed permutation matrices. A direct calculation shows that the resulting matrix is diagonal, and its diagonal entries are identically equal to $\bar{\text{tr}}[\mathbb{E} f(\mathbf{X})]$. \square

We also require a trace inequality that is related to the mean value theorem. This result specializes [MJC⁺12, Lem. 3.4].

Proposition 7.2 (Mean Value Trace Inequality). *Let $f : I \rightarrow \mathbb{R}$ be a function on an interval I of the real line whose derivative f' is convex. For all $\mathbf{A}, \mathbf{B} \in \mathbb{H}^d(I)$,*

$$\bar{\text{tr}}[(\mathbf{A} - \mathbf{B})(f(\mathbf{A}) - f(\mathbf{B}))] \leq \frac{1}{2} \bar{\text{tr}}[(\mathbf{A} - \mathbf{B})^2 \cdot (f'(\mathbf{A}) + f'(\mathbf{B}))].$$

7.2 Proof of Theorem 2.8

The argument proceeds in three steps. First, we present some elements of the matrix Laplace transform method. Second, we use the subadditivity of matrix φ -entropy to deduce a differential inequality for the trace moment generating function of the random matrix. Finally, we explain how to integrate the differential inequality to obtain the concentration result.

7.2.1 The Matrix Laplace Transform Method

We begin with a matrix extension of the moment generating function (mgf), which has played a major role in recent work on matrix concentration.

Definition 7.3 (Trace Mgf). *Let \mathbf{Y} be a random Hermitian matrix. The normalized trace moment generating function of \mathbf{Y} is defined as*

$$m(\theta) := m_{\mathbf{Y}}(\theta) := \mathbb{E} \bar{\text{tr}} e^{\theta \mathbf{Y}} \quad \text{for } \theta \in \mathbb{R}.$$

The expectation need not exist for all values of θ .

The following proposition explains how the trace mgf can be used to study the maximum eigenvalue of a random Hermitian matrix [Tro11, Prop. 3.1].

Proposition 7.4 (Matrix Laplace Transform Method). *Let $\mathbf{Y} \in \mathbb{H}^d$ be a random matrix with normalized trace mgf $m(\theta) := \bar{\text{tr}} e^{\theta \mathbf{Y}}$. For each $t \in \mathbb{R}$,*

$$\mathbb{P} \{ \lambda_{\max}(\mathbf{Y}) \geq t \} \leq d \cdot \inf_{\theta > 0} e^{-\theta t + \log m(\theta)}.$$

7.2.2 A Differential Inequality for the Trace Mgf

Suppose that $\mathbf{Y} \in \mathbb{H}^d$ is a random Hermitian matrix that depends on a random vector $\mathbf{x} := (X_1, \dots, X_n)$. We require the distribution of \mathbf{Y} to be invariant under signed permutations, and we insist that $\|\mathbf{Y}\|$ is bounded. Without loss of generality, assume that \mathbf{Y} has zero mean. Throughout the argument, we let the notation of Section 2.2.3 and Theorem 2.8 prevail.

Let us explain how to use the subadditivity of matrix φ -entropy to derive a differential inequality for the trace mgf. Consider the function $\varphi(t) = t \log t$, which belongs to the Φ_∞ class because of Theorem 2.3(1). Introduce the random positive-definite matrix $\mathbf{Z} := e^{\theta \mathbf{Y}}$, where $\theta > 0$. We write out an expression for the matrix φ -entropy of \mathbf{Z} :

$$\begin{aligned} H_\varphi(\mathbf{Z}) &= \mathbb{E} \bar{\text{tr}}[\varphi(\mathbf{Z}) - \varphi(\mathbb{E} \mathbf{Z})] \\ &= \mathbb{E} \bar{\text{tr}}[(\theta \mathbf{Y})e^{\theta \mathbf{Y}} - e^{\theta \mathbf{Y}} \log \mathbb{E} e^{\theta \mathbf{Y}}] \\ &= \theta \cdot \mathbb{E} \bar{\text{tr}}[\mathbf{Y}e^{\theta \mathbf{Y}}] - (\mathbb{E} \bar{\text{tr}} e^{\theta \mathbf{Y}}) \log(\mathbb{E} \bar{\text{tr}} e^{\theta \mathbf{Y}}) \\ &= \theta m'(\theta) - m(\theta) \log m(\theta). \end{aligned} \tag{7.2}$$

In the third line, we have applied Lemma 7.1 to the logarithm in the second term, relying on the fact that \mathbf{Y} is invariant under signed permutations. To reach the last line, we recognize that $m'(\theta) = \mathbb{E} \bar{\text{tr}}(\mathbf{Y}e^{\theta \mathbf{Y}})$. We have used the boundedness of $\|\mathbf{Y}\|$ to justify this derivative calculation.

Corollary 2.6 provides an upper bound for the matrix φ -entropy. Define the derivative $\psi(t) = \varphi'(t) = 1 + \log t$. Then

$$\begin{aligned} H_\varphi(\mathbf{Z}) &\leq \frac{1}{2} \sum_{i=1}^n \mathbb{E} \bar{\text{tr}}[(\mathbf{Z} - \mathbf{Z}'_i)(\psi(\mathbf{Z}) - \psi(\mathbf{Z}'_i))] \\ &= \frac{\theta}{2} \sum_{i=1}^n \mathbb{E} \bar{\text{tr}}[(e^{\theta \mathbf{Y}} - e^{\theta \mathbf{Y}'_i})(\mathbf{Y} - \mathbf{Y}'_i)]. \end{aligned}$$

Consider the function $f : t \mapsto e^{\theta t}$. Its derivative $f' : t \mapsto \theta e^{\theta t}$ is convex because $\theta > 0$, so Proposition 7.2 delivers the bound

$$\begin{aligned} H_\varphi(\mathbf{Z}) &\leq \frac{\theta^2}{4} \sum_{i=1}^n \mathbb{E} \bar{\text{tr}}[(e^{\theta \mathbf{Y}} + e^{\theta \mathbf{Y}'_i})(\mathbf{Y} - \mathbf{Y}'_i)^2] \\ &= \frac{\theta^2}{2} \sum_{i=1}^n \mathbb{E} \bar{\text{tr}}[e^{\theta \mathbf{Y}}(\mathbf{Y} - \mathbf{Y}'_i)^2] \\ &= \frac{\theta^2}{2} \sum_{i=1}^n \mathbb{E} \bar{\text{tr}}[e^{\theta \mathbf{Y}} \cdot \mathbb{E}[(\mathbf{Y} - \mathbf{Y}'_i)^2 | \mathbf{x}]]. \end{aligned}$$

The second relation follows from the fact that \mathbf{Y} and \mathbf{Y}'_i are exchangeable, conditional on \mathbf{x}_{-i} . The last line is just the tower property of conditional expectation, combined with the observation that \mathbf{Y} is a function of \mathbf{x} . To continue, we simplify the expression and make some additional bounds.

$$\begin{aligned} H_\varphi(\mathbf{Z}) &\leq \frac{\theta^2}{2} \mathbb{E} \bar{\text{tr}} \left[e^{\theta \mathbf{Y}} \cdot \sum_{i=1}^n \mathbb{E}[(\mathbf{Y} - \mathbf{Y}'_i)^2 | \mathbf{x}] \right] \\ &\leq \frac{\theta^2}{2} (\mathbb{E} \bar{\text{tr}} e^{\theta \mathbf{Y}}) \left\| \sum_{i=1}^n \mathbb{E}[(\mathbf{Y} - \mathbf{Y}'_i)^2 | \mathbf{x}] \right\| \\ &\leq \frac{\theta^2 V_{\mathbf{Y}}}{2} \cdot m(\theta). \end{aligned} \tag{7.3}$$

The second relation follows from a standard trace inequality and the observation that $e^{\theta \mathbf{Y}}$ is positive definite. Last, we identify the variance measure $V_{\mathbf{Y}}$ defined in (2.4) and the trace mgf $m(\theta)$.

Combine the expression (7.2) with the inequality (7.3) to arrive at the estimate

$$\theta m'(\theta) - m(\theta) \log m(\theta) \leq \frac{\theta^2 V_{\mathbf{Y}}}{2} \cdot m(\theta) \quad \text{for } \theta > 0. \quad (7.4)$$

We can use this differential inequality to obtain bounds on the trace mgf $m(\theta)$.

7.2.3 Solving the Differential Inequality

Rearrange the differential inequality (7.4) to obtain

$$\frac{d}{d\theta} \left[\frac{\log m(\theta)}{\theta} \right] = \frac{m'(\theta)}{\theta m(\theta)} - \frac{\log m(\theta)}{\theta^2} \leq \frac{V_{\mathbf{Y}}}{2}. \quad (7.5)$$

The l'Hôpital rule allows us to calculate the value of $\theta^{-1} \log m(\theta)$ at zero. Since $m(0) = 1$,

$$\lim_{\theta \rightarrow 0} \frac{\log m(\theta)}{\theta} = \lim_{\theta \rightarrow 0} \frac{m'(\theta)}{m(\theta)} = \lim_{\theta \rightarrow 0} \frac{\mathbb{E} \bar{\text{tr}}(\mathbf{Y} e^{\theta \mathbf{Y}})}{\mathbb{E} \bar{\text{tr}} e^{\theta \mathbf{Y}}} = \mathbb{E} \bar{\text{tr}} \mathbf{Y} = 0.$$

This is where we use the hypothesis that \mathbf{Y} has mean zero. Now, we integrate (7.5) from zero to some positive value θ to find that the trace mgf satisfies

$$\frac{\log m(\theta)}{\theta} \leq \frac{\theta V_{\mathbf{Y}}}{2} \quad \text{when } \theta > 0. \quad (7.6)$$

The approach in this section is usually referred to as the Herbst argument [Led99].

7.2.4 The Laplace Transform Argument

We are now prepared to finish the argument. Combine the matrix Laplace transform method, Proposition 7.4, with the trace mgf bound (7.6) to reach

$$\mathbb{P} \{ \lambda_{\max}(\mathbf{Y}) \geq t \} \leq d \cdot \inf_{\theta > 0} e^{-\theta t + \log m(\theta)} \leq d \cdot \inf_{\theta > 0} e^{-\theta t + \theta^2 V_{\mathbf{Y}}/2} = d \cdot e^{-t^2/(2V_{\mathbf{Y}})}. \quad (7.7)$$

To obtain the result for the minimum eigenvalue, we note that

$$\mathbb{P} \{ \lambda_{\min}(\mathbf{Y}) \leq -t \} = \mathbb{P} \{ \lambda_{\max}(-\mathbf{Y}) \geq t \} \leq d \cdot e^{-t^2/(2V_{\mathbf{Y}})}.$$

The inequality follows when we apply (7.7) to the random matrix $-\mathbf{Y}$. This completes the proof of Theorem 2.8.

8 Moment Inequalities for Random Matrices with Bounded Differences

In this section, we prove Theorem 2.9, which gives information about the moments of a random matrix that satisfies a kind of self-bounding property.

Proof of Theorem 2.9. Fix a number $q \in \{2, 3, 4, \dots\}$. Suppose that $\mathbf{Y} \in \mathbb{H}_+^d$ is a random positive-semidefinite matrix that depends on a random vector $x := (X_1, \dots, X_n)$. We require the distribution of \mathbf{Y} to be invariant under signed permutations, and we assume that $\mathbb{E}(\|\mathbf{Y}\|^q) < \infty$. The notation of Section 2.2.3 and Theorem 2.9 remains in force.

Let us explain how the subadditivity of matrix φ -entropy leads to a bound on the q th trace moment of \mathbf{Y} . Consider the power function $\varphi(t) = t^{q/(q-1)}$. Theorem 6.7 ensures

that $\varphi \in \Phi_\infty$ because $q/(q-1) \in (1, 2]$. Introduce the random positive-semidefinite matrix $\mathbf{Z} := \mathbf{Y}^{q-1}$. Then

$$\begin{aligned} H_\varphi(\mathbf{Z}) &= \mathbb{E} \bar{\text{tr}} [\varphi(\mathbf{Z}) - \varphi(\mathbb{E} \mathbf{Z})] \\ &= \mathbb{E} \bar{\text{tr}}(\mathbf{Y}^q) - \bar{\text{tr}} [(\mathbb{E}(\mathbf{Y}^{q-1}))^{q/(q-1)}] \\ &= \mathbb{E} \bar{\text{tr}}(\mathbf{Y}^q) - [\mathbb{E} \bar{\text{tr}}(\mathbf{Y}^{q-1})]^{q/(q-1)}. \end{aligned} \tag{8.1}$$

The transition to the last line requires Lemma 7.1.

Corollary 2.6 provides an upper bound for the matrix φ -entropy. Define the derivative $\psi(t) = \varphi'(t) = (q/(q-1)) \cdot t^{1/(q-1)}$. We have

$$\begin{aligned} H_\varphi(\mathbf{Z}) &\leq \frac{1}{2} \sum_{i=1}^n \mathbb{E} \bar{\text{tr}} [(\mathbf{Z} - \mathbf{Z}'_i)(\psi(\mathbf{Z}) - \psi(\mathbf{Z}'_i))] \\ &= \frac{q}{2(q-1)} \sum_{i=1}^n \mathbb{E} \bar{\text{tr}} [(\mathbf{Y}^{q-1} - (\mathbf{Y}'_i)^{q-1})(\mathbf{Y} - \mathbf{Y}'_i)] \end{aligned}$$

The function $f : t \mapsto t^{q-1}$ has the derivative $f' : t \mapsto (q-1)t^{q-2}$, which is convex because $q \in \{2, 3, 4, \dots\}$. Therefore, the mean value trace inequality, Proposition 7.2, delivers the bound

$$\begin{aligned} H_\varphi(\mathbf{Z}) &\leq \frac{q}{4} \sum_{i=1}^n \mathbb{E} \bar{\text{tr}} [(\mathbf{Y}^{q-2} + (\mathbf{Y}'_i)^{q-2})(\mathbf{Y} - \mathbf{Y}'_i)^2] \\ &= \frac{q}{2} \sum_{i=1}^n \mathbb{E} \bar{\text{tr}} [\mathbf{Y}^{q-2}(\mathbf{Y} - \mathbf{Y}'_i)^2] \\ &= \frac{q}{2} \sum_{i=1}^n \mathbb{E} \bar{\text{tr}} [\mathbf{Y}^{q-2} \mathbb{E}[(\mathbf{Y} - \mathbf{Y}'_i)^2 | \mathbf{x}]]. \end{aligned}$$

The second identity holds because \mathbf{Y} and \mathbf{Y}'_i are exchangeable, conditional on \mathbf{x}_{-i} . The last line follows from the tower property of conditional expectation. We simplify this expression as follows.

$$\begin{aligned} H_\varphi(\mathbf{Z}) &\leq \frac{q}{2} \mathbb{E} \bar{\text{tr}} \left[\mathbf{Y}^{q-2} \cdot \sum_{i=1}^n \mathbb{E}[(\mathbf{Y} - \mathbf{Y}'_i)^2 | \mathbf{x}] \right] \\ &\leq \frac{q}{2} \mathbb{E} \bar{\text{tr}} [\mathbf{Y}^{q-2} \cdot c\mathbf{Y}] \\ &= \frac{cq}{2} \mathbb{E} \bar{\text{tr}}(\mathbf{Y}^{q-1}). \end{aligned} \tag{8.2}$$

The second inequality derives from the hypothesis (2.5) that $\mathbf{V}_\mathbf{Y} \preceq c\mathbf{Y}$. Note that this bound requires the fact that \mathbf{Y}^{q-2} is positive semidefinite.

Combine the expression (8.1) for the matrix φ -entropy with the upper bound (8.2) to achieve the estimate

$$\mathbb{E} \bar{\text{tr}}(\mathbf{Y}^q) - [\mathbb{E} \bar{\text{tr}}(\mathbf{Y}^{q-1})]^{q/(q-1)} \leq \frac{cq}{2} \mathbb{E} \bar{\text{tr}}(\mathbf{Y}^{q-1}).$$

Rewrite this bound, and invoke the numerical fact $1 + aq \leq (1 + a)^q$ to obtain

$$\begin{aligned} \mathbb{E} \bar{\text{tr}}(\mathbf{Y}^q) &\leq [\mathbb{E} \bar{\text{tr}}(\mathbf{Y}^{q-1})]^{q/(q-1)} \left(1 + \frac{cq/2}{[\mathbb{E} \bar{\text{tr}}(\mathbf{Y}^{q-1})]^{1/(q-1)}} \right) \\ &\leq [\mathbb{E} \bar{\text{tr}}(\mathbf{Y}^{q-1})]^{q/(q-1)} \left(1 + \frac{c/2}{[\mathbb{E} \bar{\text{tr}}(\mathbf{Y}^{q-1})]^{1/(q-1)}} \right)^q. \end{aligned}$$

Extract the q th root from both sides to reach

$$[\mathbb{E} \bar{\text{tr}}(\mathbf{Y}^q)]^{1/q} \leq [\mathbb{E} \bar{\text{tr}}(\mathbf{Y}^{q-1})]^{1/(q-1)} + \frac{c}{2}.$$

We have compared the q th trace moment of \mathbf{Y} with the $(q - 1)$ th trace moment. Proceeding by iteration, we arrive at

$$[\mathbb{E} \bar{\text{tr}}(\mathbf{Y}^q)]^{1/q} \leq \mathbb{E} \bar{\text{tr}} \mathbf{Y} + \frac{q-1}{2} \cdot c.$$

This observation completes the proof of Theorem 2.9. □

A Lemma 4.1, The General Case

In this appendix, we explain how to prove Lemma 4.1 in full generality. The argument calls for a simple but powerful result, known as the generalized Klein inequality [Pet94, Prop. 3], which allows us to lift a large class of scalar inequalities to matrices.

Proposition A.1 (Generalized Klein Inequality). *For each $k = 1, \dots, n$, suppose that $f_k : I_1 \rightarrow \mathbb{R}$ and $g_k : I_2 \rightarrow \mathbb{R}$ are functions on intervals I_1 and I_2 of the real line. Suppose that*

$$\sum_{k=1}^n f_k(a) g_k(b) \geq 0 \quad \text{for all } a \in I_1 \text{ and } b \in I_2.$$

Then, for each natural number d ,

$$\sum_{k=1}^n \bar{\text{tr}}[f_k(\mathbf{A}) g_k(\mathbf{B})] \geq 0 \quad \text{for all } \mathbf{A} \in \mathbb{H}^d(I_1) \text{ and } \mathbf{B} \in \mathbb{H}^d(I_2).$$

Proof of Lemma 4.1, General Case. We retain the notation from Lemma 4.1. In particular, we assume that \mathbf{Z} is a random positive-definite matrix for which $\|\mathbf{Z}\|$ and $\|\varphi(\mathbf{Z})\|$ are both integrable. We also assume that \mathbf{T} is a random positive-definite matrix with $\|\mathbf{T}\|$ and $\|\varphi(\mathbf{T})\|$ integrable.

For $n \in \mathbb{N}$, define the function $l_n(a) := (a \vee 1/n) \wedge n$, where \vee denotes the maximum operator and \wedge denotes the minimum operator. Consider the random matrices $\mathbf{Z}_n := l_n(\mathbf{Z})$ and $\mathbf{T}_k := l_k(\mathbf{T})$ for each $k, n \in \mathbb{N}$. These matrices have eigenvalues that are bounded and bounded away from zero, so these entities satisfy the inequality (4.3) we have already established.

$$H_\varphi(\mathbf{Z}_n) \geq \mathbb{E} \bar{\text{tr}} [(\psi(\mathbf{T}_k) - \psi(\mathbb{E} \mathbf{T}_k))(\mathbf{Z}_n - \mathbf{T}_k) + \mathbb{E} \varphi(\mathbf{T}_k - \varphi(\mathbb{E} \mathbf{T}_k))].$$

Rearrange the terms in this inequality to obtain

$$\mathbb{E} \bar{\text{tr}} \Gamma(\mathbf{Z}_n, \mathbf{T}_k) \geq \bar{\text{tr}} [-\psi(\mathbb{E} \mathbf{T}_k)(\mathbb{E} \mathbf{Z}_n - \mathbb{E} \mathbf{T}_k) - \varphi(\mathbb{E} \mathbf{T}_k) + \varphi(\mathbb{E} \mathbf{Z}_n)], \tag{A.1}$$

where we have introduced the function

$$\Gamma(\mathbf{A}, \mathbf{B}) := \varphi(\mathbf{A}) - \varphi(\mathbf{B}) - (\mathbf{A} - \mathbf{B})\psi(\mathbf{B}) \quad \text{for } \mathbf{A}, \mathbf{B} \in \mathbb{H}_{++}^d.$$

To complete the proof of Lemma 4.1, we must develop the bound

$$\mathbb{E} \bar{\text{tr}} \Gamma(\mathbf{Z}, \mathbf{T}) \geq \bar{\text{tr}} [-\psi(\mathbb{E} \mathbf{T})(\mathbb{E} \mathbf{Z} - \mathbb{E} \mathbf{T}) - \varphi(\mathbb{E} \mathbf{T}) + \varphi(\mathbb{E} \mathbf{Z})] \tag{A.2}$$

by driving $k, n \rightarrow \infty$ in (A.1).

Let us begin with the right-hand side of (A.1). We have the sure limit $\mathbf{Z}_n \rightarrow \mathbf{Z}$. Therefore, the Dominated Convergence Theorem guarantees that $\mathbb{E} \mathbf{Z}_n \rightarrow \mathbb{E} \mathbf{Z}$ because $\|\mathbf{Z}\|$ is integrable and $\|\mathbf{Z}_n\| \leq \|\mathbf{Z}\|$. Likewise, $\mathbb{E} \mathbf{T}_k \rightarrow \mathbb{E} \mathbf{T}$. The functions φ and ψ are continuous, so the limit of the right-hand side of (A.1) satisfies

$$\begin{aligned} & \bar{\text{tr}} [-\psi(\mathbb{E} \mathbf{T}_k)(\mathbb{E} \mathbf{Z}_n - \mathbb{E} \mathbf{T}_k) - \varphi(\mathbb{E} \mathbf{T}_k) + \varphi(\mathbb{E} \mathbf{Z}_n)] \\ & \rightarrow \bar{\text{tr}} [-\psi(\mathbb{E} \mathbf{T})(\mathbb{E} \mathbf{Z} - \mathbb{E} \mathbf{T}) - \varphi(\mathbb{E} \mathbf{T}) + \varphi(\mathbb{E} \mathbf{Z})]. \end{aligned} \tag{A.3}$$

This expression coincides with the right-hand side of (A.2).

Taking the limit of the left-hand side of (A.1) is more involved because the function ψ may grow quickly at zero and infinity. We accomplish our goal in two steps. First, we take the limit as $n \rightarrow \infty$. Afterward, we take the limit as $k \rightarrow \infty$.

Introduce the nonnegative function

$$\gamma(z, t) := \varphi(z) - \varphi(t) - (z - t)\psi(t) \quad \text{for } z, t > 0.$$

Boucheron et al. [BBLM05, p. 525] establish that

$$\gamma(l_n(z), l_k(t)) \leq \gamma(1, l_k(t)) + \gamma(z, l_k(t)) \quad \text{for } z, t > 0. \quad (\text{A.4})$$

The generalized Klein inequality, Proposition A.1, can be applied (with due diligence) to extend (A.4) to matrices. In particular,

$$\bar{\text{tr}} \Gamma(\mathbf{Z}_n, \mathbf{T}_k) = \bar{\text{tr}} \Gamma(l_n(\mathbf{Z}), l_k(\mathbf{T})) \leq \bar{\text{tr}}[\Gamma(\mathbf{I}, l_k(\mathbf{T})) + \Gamma(\mathbf{Z}, l_k(\mathbf{T}))] = \bar{\text{tr}}[\Gamma(\mathbf{I}, \mathbf{T}_k) + \Gamma(\mathbf{Z}, \mathbf{T}_k)].$$

Observe that the right-hand side of this inequality is integrable. Indeed, all of the quantities involving \mathbf{T}_k are uniformly bounded because the eigenvalues of \mathbf{T}_k fall in the range $[k^{-1}, k]$ and the functions φ and ψ are continuous on this interval. The terms involving \mathbf{Z} may not be bounded, but they are integrable because $\|\mathbf{Z}\|$ and $\|\varphi(\mathbf{Z})\|$ are integrable. We may now apply the Dominated Convergence Theorem to take the limit:

$$\mathbb{E} \bar{\text{tr}} \Gamma(\mathbf{Z}_n, \mathbf{T}_k) \rightarrow \mathbb{E} \bar{\text{tr}} \Gamma(\mathbf{Z}, \mathbf{T}_k) \quad \text{as } n \rightarrow \infty, \quad (\text{A.5})$$

where we rely again on the sure limit $\mathbf{Z}_n \rightarrow \mathbf{Z}$ as $n \rightarrow \infty$.

Boucheron et al. also establish that

$$\gamma(z, l_k(t)) \leq \gamma(z, 1) + \gamma(z, t) \quad \text{for } z, t > 0.$$

The generalized Klein inequality, Proposition A.1, ensures that

$$\bar{\text{tr}} \Gamma(\mathbf{Z}, \mathbf{T}_k) \leq \bar{\text{tr}}[\Gamma(\mathbf{Z}, \mathbf{I}) + \Gamma(\mathbf{Z}, \mathbf{T})].$$

We may assume that the second term on the right-hand side is integrable or else the desired inequality (A.2) would be vacuous. The first term is integrable because $\|\mathbf{Z}\|$ and $\|\varphi(\mathbf{Z})\|$ are integrable. Therefore, we may apply the Dominated Convergence Theorem:

$$\mathbb{E} \bar{\text{tr}} \Gamma(\mathbf{Z}, \mathbf{T}_k) \rightarrow \mathbb{E} \bar{\text{tr}} \Gamma(\mathbf{Z}, \mathbf{T}) \quad \text{as } k \rightarrow \infty, \quad (\text{A.6})$$

where we rely again on the sure limit $\mathbf{T}_k \rightarrow \mathbf{T}$ as $k \rightarrow \infty$.

In summary, the limits (A.5) and (A.6) provide that $\mathbb{E} \bar{\text{tr}} \Gamma(\mathbf{Z}_n, \mathbf{T}_k) \rightarrow \mathbb{E} \bar{\text{tr}} \Gamma(\mathbf{Z}, \mathbf{T})$ as $k, n \rightarrow \infty$. In view of the limit (A.3), we have completed the proof of (A.2). \square

References

- [AGZ10] G. W. Anderson, A. Guionnet, and O. Zeitouni. *An introduction to random matrices*, volume 118 of *Cambridge Studies in Advanced Mathematics*. Cambridge University Press, Cambridge, 2010.
- [And79] T. Ando. Concavity of certain maps on positive definite matrices and applications to Hadamard products. *Linear Algebra Appl.*, 26:203–241, 1979. MR-0535686
- [AW02] R. Ahlswede and A. Winter. Strong converse for identification via quantum channels. *IEEE Trans. Inform. Theory*, 48(3):569–579, 2002. MR-1889969
- [BBLM05] S. Boucheron, O. Bousquet, G. Lugosi, and P. Massart. Moment inequalities for functions of independent random variables. *Ann. Probab.*, 33(2):514–560, 2005. MR-2123200

Subadditivity of matrix φ -entropy

- [BG12] C. Boutsidis and A. Gittens. Improved matrix algorithms via the subsampled randomized Hadamard transform. Available at arXiv:1204.0062, 2012. MR-3101094
- [Bha97] R. Bhatia. *Matrix Analysis*, volume 169 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, 1997. MR-1477662
- [Bha07] R. Bhatia. *Positive Definite Matrices*. Princeton Series in Applied Mathematics. Princeton University Press, Princeton, NJ, 2007. MR-2284176
- [BL98] S. G. Bobkov and M. Ledoux. On modified logarithmic Sobolev inequalities for Bernoulli and Poisson measures. *J. Funct. Anal.*, 156(2):347–365, 1998. MR-1636948
- [BLM03] S. Boucheron, G. Lugosi, and P. Massart. Concentration inequalities using the entropy method. *Ann. Probab.*, 31(3):1583–1614, 2003. MR-1989444
- [BLM09] S. Boucheron, G. Lugosi, and P. Massart. On concentration of self-bounding functions. *Electron. J. Probab.*, 14:no. 64, 1884–1899, 2009. MR-2540852
- [BLM13] S. Boucheron, G. Lugosi, and P. Massart. *Concentration inequalities*. Oxford University Press, 2013.
- [Bou02] O. Bousquet. A Bennett concentration inequality and its application to suprema of empirical processes. *C. R. Math. Acad. Sci. Paris*, 334(6):495–500, 2002. MR-1890640
- [Brè66] L. M. Brègman. Relaxation method for finding a common point of convex sets and its application to optimization problems. *Dokl. Akad. Nauk SSSR*, 171:1019–1022, 1966. MR-0210291
- [Car10] E. Carlen. Trace inequalities and quantum entropy: an introductory course. In *Entropy and the quantum*, volume 529 of *Contemp. Math.*, pages 73–140. Amer. Math. Soc., Providence, RI, 2010. MR-2681769
- [CCT12] K. Chaudhuri, F. Chung, and A. Tsiatas. Spectral clustering of graphs with general degrees in the extended planted partition model. *Journal of Machine Learning Research* 2012, pages 1–23, 2012.
- [CDL13] A. Cohen, M. A. Davenport, and D. Leviatan. On the stability and accuracy of least squares approximations. *Found. Comput. Math.*, 2013. MR-3105946
- [Cha04] D. Chafaï. Entropies, convexity, and functional inequalities: on Φ -entropies and Φ -Sobolev inequalities. *J. Math. Kyoto Univ.*, 44(2):325–363, 2004. MR-2081075
- [Cha06] D. Chafaï. Binomial-Poisson entropic inequalities and the $M/M/\infty$ queue. *ESAIM Probab. Stat.*, 10:317–339 (electronic), 2006. MR-2247924
- [Csi72] I. Csiszár. A class of measures of informativity of observation channels. *Period. Math. Hungar.*, 2:191–213, 1972. Collection of articles dedicated to the memory of Alfréd Rényi, I. MR-0335152
- [DZ11] P. Drineas and A. Zouzias. A note on element-wise matrix sparsification via a matrix-valued Bernstein inequality. *Inform. Process. Lett.*, 111:385–389, 2011. MR-2760960
- [Han13] F. Hansen. Trace functions with applications in quantum physics. *J. Stat. Phys.*, 2013. MR-3163550
- [HOZ01] W. Hebisch, R. Olkiewicz, and B. Zegarliński. On upper bound for the quantum entropy. *Linear Algebra Appl.*, 329:89–96, 2001. MR-1822224
- [HZ14] F. Hansen and Z. Zhang. Characterisation of matrix entropies. Available at arXiv:1402.2118, Feb. 2014.
- [JX03] M. Junge and Q. Xu. Noncommutative Burkholder/Rosenthal inequalities. *Ann. Probab.*, 31(2):948–995, 2003. MR-1964955
- [JX08] M. Junge and Q. Xu. Noncommutative Burkholder/Rosenthal inequalities II: Applications. *Israel J. Math.*, 167:227–282, 2008. MR-2448025
- [JZ11] M. Junge and Q. Zheng. Noncommutative Bennett and Rosenthal inequalities. Available at arXiv:1111.1027, Nov. 2011.
- [KA80] F. Kubo and T. Ando. Means of positive linear operators. *Math. Ann.*, 246(3):205–224, 1979/80. MR-0563399
- [Led99] M. Ledoux. Concentration of measure and logarithmic Sobolev inequalities. In *Séminaire de Probabilités, XXXIII*, volume 1709 of *Lecture Notes in Math.*, pages 120–216. Springer, Berlin, 1999. MR-1767995

- [Led01] M. Ledoux. *The Concentration of Measure Phenomenon*, volume 89 of *Mathematical Surveys and Monographs*. American Mathematical Society, Providence, RI, 2001. MR-1849347
- [Led97] M. Ledoux. On Talagrand's deviation inequalities for product measures. *ESAIM Probab. Statist.*, 1:63–87 (electronic), 1995/97. MR-1399224
- [Lie73] E. H. Lieb. Convex trace functions and the Wigner-Yanase-Dyson conjecture. *Advances in Math.*, 11:267–288, 1973. MR-0332080
- [Lie75] E. H. Lieb. Some convexity and subadditivity properties of entropy. *Bull. Amer. Math. Soc.*, 81:1–13, 1975. MR-0356797
- [Lin73] G. Lindblad. Entropy, information, and quantum measurements. *Commun. Math. Phys.*, 33:305–322, 1973. MR-0337240
- [LO00] R. Latała and K. Oleszkiewicz. Between Sobolev and Poincaré. In *Geometric aspects of functional analysis*, volume 1745 of *Lecture Notes in Math.*, pages 147–168. Springer, Berlin, 2000. MR-1796718
- [LP86] F. Lust-Piquard. Inégalités de Khintchine dans C_p ($1 < p < \infty$). *C. R. Acad. Sci. Paris Sér. I Math.*, 303(7):289–292, 1986. MR-0859804
- [LR73] E. H. Lieb and M. B. Ruskai. Proof of the strong subadditivity of quantum-mechanical entropy. *J. Math. Phys.*, 14(12):1938–1941, 1973. MR-0345558
- [Mas00a] P. Massart. About the constants in Talagrand's concentration inequalities for empirical processes. *Ann. Probab.*, 28(2):863–884, 2000. MR-1782276
- [Mas00b] P. Massart. Some applications of concentration inequalities to statistics. *Ann. Fac. Sci. Toulouse Math. (6)*, 9(2):245–303, 2000. Probability theory. MR-1813803
- [Min12] S. Minsker. On some extensions of Bernstein's inequality for self-adjoint operators. Available at arXiv:1112.5448, Jan. 2012.
- [MJC⁺12] L. Mackey, M. I. Jordan, R. Y. Chen, B. Farrell, and J. A. Tropp. Matrix concentration inequalities via the method of exchangeable pairs. Available at arXiv:1201.6002, 2012.
- [Oli09] R. I. Oliveira. Concentration of the adjacency matrix and of the Laplacian in random graphs with independent edges. Available at arXiv:0911.0600, 2009.
- [Pet94] D. Petz. A survey of certain trace inequalities. In *Functional analysis and operator theory (Warsaw, 1992)*, volume 30 of *Banach Center Publ.*, pages 287–298. Polish Acad. Sci., Warsaw, 1994. MR-1285615
- [PMT13] D. Paulin, L. Mackey, and J. A. Tropp. Deriving matrix concentration inequalities from kernel couplings. Available at arXiv:1305.0612, May 2013.
- [Rén61] A. Rényi. On measures of entropy and information. In *Proc. 4th Berkeley Sympos. Math. Statist. and Prob.*, Vol. I, pages 547–561. Univ. California Press, Berkeley, Calif., 1961. MR-0132570
- [Rio01] E. Rio. Inégalités de concentration pour les processus empiriques de classes de parties. *Probab. Theory Related Fields*, 119(2):163–175, 2001. MR-1818244
- [RS13] M. Raginsky and I. Sason. Concentration of measure inequalities in information theory, communications and coding. *Foundations and Trends in Communications and Information Theory*, 10(1–2):1–246, 2013. Available at arXiv:1212.4663.
- [Rud99] M. Rudelson. Random vectors in the isotropic position. *J. Funct. Anal.*, 164(1):60–72, 1999. MR-1694526
- [RW11] M. D. Reid and R. C. Williamson. Information, divergence and risk for binary experiments. *J. Mach. Learn. Res.*, 12:731–817, 2011. MR-2786911
- [Tal96] M. Talagrand. A new look at independence. *Ann. Probab.*, 24(1):1–34, 1996. MR-1387624
- [Tro11] J. A. Tropp. Freedman's inequality for matrix martingales. *Electron. Commun. Probab.*, 16:262–270, 2011. MR-2802042
- [Tro12a] J. A. Tropp. Improved analysis of the subsampled randomized Hadamard transform. *Adv. Adapt. Data Anal.*, 3(1–2):115–126, 2012. MR-2835584
- [Tro12b] J. A. Tropp. User-friendly tail bounds for sums of random matrices. *Found. Comput. Math.*, 12(4):389–434, 2012. MR-2946459

Subadditivity of matrix φ -entropy

[Tro12c] J. A. Tropp. User-friendly tools for random matrices: An introduction. Available at <http://users.cms.caltech.edu/~jtropp/notes/Tro12-User-Friendly-Tools-NIPS.pdf>, Dec. 2012.

[Voi02] D. Voiculescu. Free entropy. *Bull. London Math. Soc.*, 34(3):257–278, 2002. MR-1887698

Acknowledgments. RYC and JAT are with the Department of Computing and Mathematical Sciences, California Institute of Technology. JAT gratefully acknowledges support from ONR awards N00014-08-1-0883 and N00014-11-1002, AFOSR award FA9550-09-1-0643, and a Sloan Research Fellowship. JAT also wishes to thank the Moore Foundation.