# SHRINKAGE ESTIMATION OF FUNCTIONS OF LARGE NOISY SYMMETRIC MATRICES

PANAGIOTIS LOLAS [*1] AND LEXING YING [†1]

ABSTRACT. We study the problem of estimating functions of a large symmetric matrix $A$ when we only have access to a noisy estimate $\hat{A}_n = A_n + \sigma Z_n/\sqrt{n}$. We are interested in the case that $Z_n$ is a Wigner ensemble and suggest an algorithm based on nonlinear shrinkage of the eigenvalues of $\hat{A}_n$. As an intermediate step we explain how recovery of the spectrum of $A_n$ is possible using only the spectrum of $\hat{A}_n$. Our algorithm has important applications, for example, in solving high-dimensional noisy systems of equations or symmetric matrix denoising. Throughout our analysis we rely on tools from random matrix theory.

## 1. INTRODUCTION

1.1. **Problem and Assumptions.** Let $A_n \in \mathbb{R}^{n \times n}$ be a real symmetric matrix (deterministic or random), which is unknown. Instead, we have access to a noisy estimate $\hat{A}_n = A_n + \sigma n^{-1/2} Z_n$. We will often omit the subscript $n$ in our notation. We will denote by $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$ the eigenvalues of $A_n$, $w_1, \cdots, w_n$ the corresponding eigenvectors and the empirical spectral distribution of $A_n$ by $\mu_n$. The latter is the measure $\mu_n = n^{-1} \sum_{k=1}^{n} \delta_{\lambda_k}$. Similarly we are going to denote by $\hat{\lambda}_1 \geq \cdots \geq \hat{\lambda}_n$ the eigenvalues of $\hat{A}_n$ and $\hat{w}_1, \cdots, \hat{w}_n$ the corresponding eigenvectors.

**Assumptions 1.1.** *We assume that $A_n$ and $Z_n$ satisfy the following assumptions.*

(1) *The dimension $n$ of the matrix $A_n$ goes to infinity.*
(2) *The spectral distribution of the eigenvalues of $A_n$ converges weakly almost surely to a deterministic probability measure $H$.*
(3) *The measure $H$ is supported on a compact interval contained in $\mathbb{R}$ and eventually all of the eigenvalues of $A$ lie in a compact subset $[h_1, h_2]$ of $\mathbb{R}$.*
(4) *The matrix $Z_n$ is real symmetric and independent of $A_n$. The matrix $Z_n$ is a submatrix of an infinite matrix $(Z_{ij})_{1 \leq i,j \leq n}$ whose upper half has i.i.d. entries with mean 0, variance 1 and finite fourth moments.*

We are interested in estimating $h(A_n)$, where $h$ is a continuous function defined on an open set that contains $[h_1, h_2]$. Special cases include, for example, $h(x) = x$ (which is the problem of denoising $\hat{A}_n$), or $h(x) = x^{-1}$, which is interesting for solving noisy linear systems of equations. Other interesting choices might include $h(x) = \sqrt{x}$ (estimating the square root of a positive semi-definite matrix), or $h(x) = x/(x^2 + \lambda^2)$ (for estimating the regularized inverse of a symmetric matrix).

---
[*]panagd@stanford.edu.

[†]lexing@stanford.edu.

[1] Department of Mathematics, Stanford University.

1

1.2. **Our Contributions.** The main contributions of our paper are listed below:

(1) We derive (in closed form) the optimal nonlinear shrinkage for estimating $h(A_n)$ in Frobenius loss.

(2) We suggest a practical algorithm that asymptotically estimates the optimal nonlinear shrinkage for any choice of function $h$.

(3) We study the problem of recovering the limiting spectral distribution $H$ of the matrix $A_n$. We consider the cases of known and unknown noise level $\sigma^2$. Recovering the measure $H$ is important for the implementation of our algorithm.

(4) We show how our results can be used to derive the optimal shrinkage function with alternative choices of losses.

(5) We study asymptotic expansions of the optimal shrinkers when $\sigma \to 0$ and $\sigma \to \infty$.

1.3. **Related Work.** Shrinkage methods have been used in statistics in different settings with great success. In James and Stein [1992] the authors showed how estimation of the mean of a Gaussian distribution in more than 2 dimensions can be improved significantly by shrinkage of the sample estimates. For the purpose of covariance matrix estimation, linear shrinkage methods were used in Ledoit and Wolf [2004] to suggest a well-conditioned estimator of a high-dimensional covariance matrix. Using tools from random matrix theory, in Ledoit et al. [2012] the authors showed how nonlinear shrinkage methods can be used to greatly improve estimation and a nonparametric procedure that achieves greater speed and numerical stability was suggested in Ledoit et al. [2020]. For the case of spiked models, Donoho et al. [2018] used nonlinear shrinkage to estimate the population covariance matrix and derived the optimal shrinker for 26 losses, for most of them in closed form. For regularization of linear discriminant analysis, general nonlinear eigenvalue shrinkage was used in Lolas [2020] to improve the classification accuracy when the feature dimensionality is comparable to the number of samples and sharp classification error asymptotics for any shrinkage function were derived.

For the case of a deformed Wigner model as the one we consider here, Donoho and Gavish [2013] showed how eigenvalue shrinkage can be used for symmetric matrix denoising in the case that $A$ is low-rank. For the problem considered here, $h(x) = x$ was studied by Bun et al. [2016], where the authors derived the optimal nonlinear shrinkage in closed form using replica symmetry. In that case the authors showed that, given $\sigma$, the optimal shrinker depends on $H$ only through the Stieltjes transform of the limiting spectral distribution of $\hat{A}_n$. This phenomenon makes the optimal shrinkage function easy to estimate (for example, with a similar nonparametric procedure as in Ledoit et al. [2020]).

The problem of numerical computation of the free-convolution of two probability measures has been studied in Rao and Edelman [2008], Olver and Nadakuditi [2012]. The inverse problem, namely spectrum recovery (which we study for the deformed Wigner case in Section 4), has been well-studied for covariance matrices. In El Karoui et al. [2008] a convex optimization approach was used to recover population spectra from samples. In Kong et al. [2017], the authors used a moment method that works even in the sublinear regime where the dimension of the covariance matrix is much larger than the number of samples. Ledoit and Wolf [2015] used an approach that exploits the natural discreteness of population spectra and suggested solving a nonlinear optimization problem which essentially matches the

empirical eigenvalues to the quantiles of the Marcenko-Pastur distribution. The idea of natural discreteness of the population spectrum will also be useful for the case of additive free-convolution with a semicircular distribution that we consider here.

Finally, from a Bayesian perspective shrinkage methods have been considered in other settings. In a closely related problem in Etter and Ying [2020] the authors suggested a Bayesian shrinkage method to solve noisy elliptic systems of equations. For the case of covariance matrix estimation, linear shrinkage is motivated by imposing am inverse Wishart prior, while other more sophisticated priors give rise to nonlinear shrinkage methods (Yang and Berger [1994],Berger et al. [2020]).

1.4. **Organization of the Paper.** In Section 2 we review some well-known results from random matrix theory and present a new result about trace functionals that involve both $A_n$ and $\hat{A}_n$. These are going to be the essential tools that we will need for the rest of the paper. In Section 3 we derive the oracle nonlinear shrinkage estimators for general continuous functions of $A$ and asymptotic equivalents that are amenable to estimation. We also suggest an algorithm to perform asymptotically optimal nonlinear shrinkage, when $H, \sigma$ are known. Section 4 considers the problem of recovering $H, \sigma$. Firstly, we show how $H$ can be recovered, given $\sigma$, using a nonlinear optimization problem and provide theoretical guarantees for consistency. We then explain how $\sigma$ can be consistently estimated for a class of probability measures $H$. In Section 5 we study asymptotic expansions of the shrinkers and the losses when $\sigma \to 0$ and $\sigma \to \infty$. Simulations and numerical experiments are presented in Section 6. Finally, Section 7 presents the complete proofs of our results.

## 2. Almost Sure Limits for a Class of Trace Functionals

In this section we present some useful tools from random matrix theory. We start by introducing our notation and stating well-known theorems. After that, we provide some new results about asymptotics of trace functionals that include both $A$ and $\hat{A}$ which will be essential for justifying the main algorithm in Section 3.

For a probability measure $\mu$ supported on the real line we will denote its Stieltjes transform by $m_\mu(z) = \int (x - z)^{-1} \mu(dx), z \in \mathbb{C}^+$. We will often omit the measure from the subscript and just write $m(z)$, provided that it is clear which measure we are referring to. We have the following well-known result, the so-called Wigner semicircle law (Wigner [1958]).

**Theorem 2.1** (Theorem 2.4.2 in Tao [2012]). *Let* $(M_{ij})_{1 \leq i,j}$ *be mean 0, variance 1 real random variables such that* $M_{ij} = M_{ji}$ *and* $(M_{ij})_{i<j}$ *are independent and identically distributed. Then, the spectral distribution of the sequence of random matrices* $n^{-1/2} M_n = (n^{-1/2} M_{ij})_{1 \leq i,j \leq n}$ *converges weakly almost surely to the Wigner semicircular distribution:*

$$\mu_{sc} = \frac{\sqrt{(4 - x^2)_+}}{2\pi} dx.$$

The above result gives the limiting spectral distribution of Wigner matrices. For the case of a deformed Wigner matrix, such as $\hat{A}_n = A_n + \sigma n^{-1/2} Z_n$, we have under the Assumptions 1.1 in Subsection 1.1:

**Proposition 2.1.** *The matrix $\hat{A}_n$ has a limiting spectral distribution $\hat{\mu}$, which is a deterministic probability measure with Stieltjes transform $m_{\hat{\mu}}(z)$ that satisfies:*

$$m_{\hat{\mu}}(z) = \int \frac{dH(t)}{t - z - \sigma^2 m_{\hat{\mu}}(z)}.$$

This is the formula that describes the free additive convolution $H \boxplus \rho_{sc;\sigma^2}$ of a measure with a semicircular distribution (Biane [1997]). If $H = \delta_0$, we can solve for the Stieltjes transform $m_{\hat{\mu}}(z)$ in closed form and then use Stieltjes inversion to recover the Wigner law.

The first main contribution of this paper is to extend this result in the following theorem, which is analogous to the results in Ledoit and Péché [2011] for the case of covariance matrices. As in the case of covariance matrices, when Ledoit et al. [2012] used it to estimate a covariance matrix using nonlinear shrinkage, this is going to be the main tool for theoretically justifying our algorithms. In Bun et al. [2016] a similar calculation is done using using replica symmetry for matrices corrupted by orthogonally invariant noise.

**Theorem 2.2.** *With the same assumptions as in Section 1 we have for any $z \in \mathbb{C}^+$*

$$\frac{tr\left(h(A)\left(\hat{A} - z\right)^{-1}\right)}{n} \xrightarrow{a.s.} \int \frac{h(t)dH(t)}{t - z - \sigma^2 m_{\hat{\mu}}(z)}.$$

*Here, $m_{\hat{\mu}}(z)$ is the Stieltjes transform of the free additive convolution of $H$ with a semicircular distribution of variance $\sigma^2$, as in Proposition 2.1.*

Although the theorem above was stated for a function $h$ that is continuous, it can be extended to cases with finitely many discontinuities which are not on atoms of the measure $H$. In that case, taking $h(t) = \mathbb{I}_{[a,b]}$ gives the asymptotic overlap of the eigenvectors of $A, \hat{A}$, which the authors in Bun et al. [2016] derived.

## 3. Main Results

In this section, we motivate and present the main algorithm of the paper. We start by deriving an oracle estimator that optimally approximates $h(A)$ among all rotationally invariant estimators. We also find the optimal shrinker in closed form using the results from Section 2. After that, we explain how universality, namely the fact that in the large $n$ limit the distribution of the noise does not affect the asymptotics we are interested in, can be used to simulate approximately the optimally shrunk eigenvalues.

3.1. **Optimal Rotation Invariant Estimator.** We consider the spectral decomposition of $\hat{A}$, which has eigenvalues $\hat{\lambda}_1 \geq \cdots \geq \hat{\lambda}_n$ :

$$\hat{A} = \hat{W}\hat{\Lambda}\hat{W}^\top.$$

For a continuous function $h$, an estimator $\Psi(\hat{A})$ of $h(A)$ is rotationally invariant if $\Psi(O\hat{A}O^\top) = O\Psi(\hat{A})O^\top$ for any $n \times n$ orthogonal matrix $O$. Searching for a rotationally invariant estimator of $h(A)$ seems reasonable, if we do not have any prior information about the eigenstructure of $A$. If such information was available, we might be able to exploit it by approaching the problem in a Bayesian way. With that in mind, it also seems reasonable to consider $\Psi(\hat{A})$ with the same eigenvectors

as $\hat{A}$, such that $\Psi(\hat{A}) = \hat{W}D\hat{W}^\top$. We are interested in choosing $\Psi$ to minimize the Frobenius loss $\left\| \Psi(\hat{A}) - h(A) \right\|_F^2$. We observe that

$$\left\| \Psi(\hat{A}) - h(A) \right\|_F^2 = \left\| D - \hat{W}^\top h(A)\hat{W} \right\|_F^2,$$

which is minimized when

$$D^{(h)} = (d_1^{(h)}, \cdots, d_n^{(h)}) = \operatorname{diag}(\hat{W}^\top h(A)\hat{W}). \tag{1}$$

These clearly depends on the unknown matrix $A$ and is not straightforward to estimate from the data. In the case $h(x) = x$ the authors in Bun et al. [2016], Potters and Bouchaud [2020] show that the oracle quantities can be asymptotically approximated by deterministic quantities that depend only on the limiting spectral distribution of $\hat{A}$ and the noise $\sigma$. The authors call this remarkable phenomenon the large dimension miracle. It makes the oracle quantities amenable to estimation, for example via kernel estimation. However, such a miracle does not seem very likely in the case of a general $h$ (and it is not entirely clear how to extend to the case of unknown $\sigma$). For example, already for $h(x) = x^{-1}$, we will see that the optimal shrinkage is given by

$$f_{1/t}^*(x) \equiv \frac{x + \sigma^2 \int t^{-1} dH(t)}{(x + \sigma^2 u(x))^2 + \sigma^4 v(x)^2},$$

where $u(x)$ and $v(x)$ are the real and imaginary parts of $\lim_{\eta \to 0^+} m_{\hat{\mu}}(x + i\eta)$, i.e., $u(x) + iv(x) = \lim_{\eta \to 0^+} m_{\hat{\mu}}(x + i\eta)$. This already requires estimating $m_H(0) = \int t^{-1} dH(t)$ and it is not hard to see that for other functions the situation can get even more complicated.

**Definition 3.1.** *For a continuous function $h$ on an open interval that contains the support of $H$, we define the functions $u_h, v_h : \mathbb{R} \to \mathbb{R}$ by*

$$u_h(x) + iv_h(x) = \lim_{\eta \downarrow 0} \int \frac{h(t)dH(t)}{t - x - i\eta - \sigma^2 m_{\hat{\mu}}(x + i\eta)}.$$

**Remark 3.1.** *The limit above exists because $\lim_{\eta \downarrow 0} m_{\hat{\mu}}(x + i\eta)$ exists (Biane [1997]).*

Below derive the optimal shrinker for a general continuous function $h$.

**Theorem 3.1.** *Among all bounded continuous functions $f$ on an open interval containing $supp(H \boxplus \rho_{sc;\sigma^2})$ and the eigenvalues of $\hat{A}_n$, the minimizer $f_h^*(x)$ of the asymptotic quantity*

$$\lim_{n \to \infty} n^{-1} \left\| f(\hat{A}_n) - h(A_n) \right\|_F^2$$

*is given by $f_h^*(x) = v_h(x)/v(x)$ for $x \in supp(H \boxplus \rho_{sc;\sigma^2})$.*

According to Theorem 2.2, the measure $n^{-1} \sum_{i=1}^n d_i^{(h)} \delta_{\hat{\lambda}_i}$ converges weakly almost surely to a measure with density $\pi^{-1} v_h$. This suggests that the asymptotic analog of the oracle quantities $d_i^{(h)}$ is the quantity $v_h/v$ derived above. As an immediate corollary of Theorem 3.1 we have the following:

**Corollary 3.1.** *(1) For the choice $h(t) = t$,*

$$u_t(x) + iv_t(x) = \lim_{z = x + i\eta, \eta \downarrow 0} 1 + zm_{\hat{\mu}}(z) + \sigma^2 m_{\hat{\mu}}(z)^2. \tag{2}$$

*This gives the optimal shrinkage function $f_t^*(x) = x + 2\sigma^2 u(x)$.*

(2) *For the choice* $h(t) = 1/t$,

$$u_{1/t}(x) + iv_{1/t}(x) = \lim_{z=x+i\eta,\eta\downarrow 0} \frac{m_{\hat{\mu}}(z) - m_H(0)}{z + \sigma^2 m_{\hat{\mu}}(z)}. \tag{3}$$

*This gives the optimal shrinkage function*

$$f_{1/t}^*(x) = \frac{x + \sigma^2 m_H(0)}{(x + \sigma^2 u)^2 + \sigma^4 v(x)^2}.$$

(3) *For the choice* $h(t) = t^2$,

$$u_{t^2}(x) + iv_{t^2}(x) = \lim_{z=x+i\eta,\eta\downarrow 0} \int \frac{t^2 dH(t)}{t - z - \sigma^2 m_{\hat{\mu}}(z)}$$

$$= \lim_{z=x+i\eta,\eta\downarrow 0} \int t dH(t) + z + \sigma^2 m_{\hat{\mu}}(z) + m_{\hat{\mu}}(z)\left[z + \sigma^2 m_{\hat{\mu}}(z)\right]^2 \tag{4}$$

$$= \int t dH(t) + x + \sigma^2(u(x) + iv(x)) + (u(x) + iv(x))\left[x + \sigma^2(u(x) + iv(x))\right]^2.$$

*This gives the optimal shrinkage function*

$$f_{t^2}^*(x) = \sigma^2 + (x + \sigma^2 u(x))^2 - \sigma^4 v^2(x) + 2\sigma^2 u(x)(x + \sigma^2 u(x)).$$

**Remark 3.2.** *Using Theorem 3.1 we can show that for estimating $A_n^k$ in Frobenius norm we need the first $(k-2)$ moments of the measure $H$ for $k \geq 3$.*

3.1.1. *Pseudoinverses and Regularized Pseudoinverses.* We study the optimal shrinkage to estimate $A\left(A^2 + \lambda^2 I_n\right)^{-1}$. If $\lambda \downarrow 0$, this converges to the pseudoinverse of the matrix $A$. Using our usual notation we have $h = h(t; \lambda) = t/(t^2 + \lambda^2)$. This gives

$$\int \frac{h(t)dH(t)}{t - z - \sigma^2 m_{\hat{\mu}}(z)} = \int \frac{t dH(t)}{(t^2 + \lambda^2)(t - z - \sigma^2 m_{\hat{\mu}}(z))}$$

$$= \int \left[\frac{z + \sigma^2 m_{\hat{\mu}}(z)}{(z + \sigma^2 m_{\hat{\mu}}(z))^2 + \lambda^2}\frac{1}{t - z - \sigma^2 m_{\hat{\mu}}(z)}\right]dH(t)$$

$$+ \frac{1}{2(\lambda i - z - \sigma^2 m_{\hat{\mu}}(z))}\int \frac{dH(t)}{t - \lambda i} - \frac{1}{2(\lambda i + z + \sigma^2 m_{\hat{\mu}}(z))}\int \frac{dH(t)}{t + \lambda i} \tag{5}$$

$$= \frac{m_{\hat{\mu}}(z)(z + \sigma^2 m_{\hat{\mu}}(z))}{(z + \sigma^2 m(z))^2 + \lambda^2} + \frac{m_H(\lambda i)}{2(\lambda i - z - \sigma^2 m_{\hat{\mu}}(z))} - \frac{m_H(-\lambda i)}{2(\lambda i + z + \sigma^2 m_{\hat{\mu}}(z))}$$

which allows us to compute the optimal shrinkage as a function of $u, v, m_H(\lambda i)$.

For the case of the pseudoinverse of a Hermitian matrix $A$, we examine the following scenario. We assume that there exist fixed $\delta > 0, p \in (0, 1)$ such that $A$ has $a_n$ eigenvalues equal to 0, $n - a_n$ eigenvalues greater than $\delta$ and $a_n/n \xrightarrow{a.s.} p$ as $n \to \infty$. In that case we can write $H = p\delta_0 + (1-p)\nu$, where $\nu$ is a probability measure with support contained in $[\delta, \infty)$. Under these assumptions the pseudoinverse of $A$ can be written as a function $h(A)$, where $h$ is continuous on $[0, \infty]$, $h(x) = 1/x$ for $x \geq \delta$ and $h(x) = 0$ in an open set containing 0. We find in this case the Stieljes transform of $H \boxplus \rho_{sc;\sigma^2}$ satisfies

$$m_{\hat{\mu}}(z) = -\frac{p}{z + \sigma^2 m_{\hat{\mu}}(z)} + (1-p)\int \frac{d\nu(t)}{t - z - \sigma^2 m_{\hat{\mu}}(z)}. \tag{6}$$

Using this we see that

$$\int \frac{h(t)dH(t)}{t - z - \sigma^2 m_{\hat{\mu}}(z)}$$

$$= (1-p) \int \frac{d\nu(t)}{t(t - z - \sigma^2 m_{\hat{\mu}}(z))} = \frac{1-p}{z + \sigma^2 m_{\hat{\mu}}(z)} \int \left[ \frac{1}{t - z - \sigma^2 m_{\hat{\mu}}(z)} - \frac{1}{t} \right] d\nu(t)$$

$$= \frac{1-p}{z + \sigma^2 m_{\hat{\mu}}(z)} \left[ \frac{m_{\hat{\mu}}(z) + \frac{p}{z + \sigma^2 m_{\hat{\mu}}(z)}}{1 - p} - m_{\nu}(0) \right]$$

$$= \frac{m_{\hat{\mu}}(z)}{z + \sigma^2 m_{\hat{\mu}}(z)} + \frac{p}{(z + \sigma^2 m_{\hat{\mu}}(z))^2} - \frac{(1-p)m_{\nu}(0)}{z + \sigma^2 m_{\hat{\mu}}(z)}.$$

$$(7)$$

**Remark 3.3.** *(1) If $p = 0$, the last formula reduces to*

$$\frac{m_{\hat{\mu}}(z) - m_{\nu}(0)}{z + \sigma^2 m_{\hat{\mu}}(z)} = \frac{m_{\hat{\mu}}(z) - m_H(0)}{z + \sigma^2 m_{\hat{\mu}}(z)}.$$

*This, as expected, agrees with Corolllay 3.1.*

*(2) When $H = p\delta_0 + (1-p)\nu$, we have*

$$m_H(\lambda i) = -\frac{p}{\lambda i} + (1-p)m_{\nu}(\lambda i)$$

*and*

$$m_H(-\lambda i) = \frac{p}{\lambda i} + (1-p)m_{\nu}(-\lambda i).$$

*Using these it is straightforward to see that the optimal shrinkage for the regularized pseudoinverse converges to the optimal shrinkage for the pseudoinverse as $\lambda \downarrow 0$.*

3.2. **Monte-Carlo Nonlinear Shrinkage.** We are now going to present an algorithm to approximate the oracle quantities. Based on Theorem 3.1, it is natural to try to compute $H$ and then solve for $u_h, v_h$. Our algorithm does not require solving numerically the equation for the Stieltjes transform of the additive free convolution of $H$ with a semicircular distribution, which can be tricky (Olver and Nadakuditi [2012]). We think that the general idea behind it is likely to be applied in more complicated cases, in particular in problems that do not have simple formulas for the optimal shrinkage as derived in Theorem 3.1. The key observation is that the asymptotic equivalents of the oracle quantities only depend on $H$ and are universal for all noise distributions. Hence, although $A, Z$ are unknown, it is possible to replicate the asymptotic equivalents to the oracles using a Monte-Carlo simulation.

Suppose that we know $\sigma, \lambda_1, \cdots, \lambda_n$, or estimates $\tilde{\sigma}, \tilde{\lambda}_1, \cdots, \tilde{\lambda}_n$ of those are available. The topic of finding suitable choices for $\tilde{\sigma}$ and $\tilde{\lambda}_i$ for $1 \leq i \leq n$ is going to be the topic of the next section, as suggested by Theorem 4.1. Then, we suggest the following simple procedure in Algorithm 3.1 for approximately optimal nonlinear shrinkage of the eigenvalues of $\hat{A}$ to estimate $h(A)$ in Frobenius norm. The complexity of the algorithm is $\mathcal{O}(Kn^3)$. Notice that we use the notation $GOE(n)$ for the Gaussian Orthogonal Ensemble in $\mathbb{R}^{n \times n}$ (Tao [2012]).

Algorithm 3.1 approximates the oracle nonlinear shrinkage in the following sense.

**Theorem 3.2.** *For a bounded continuous function $h$ defined on an open set that contains the support of $H$, let $d_i^{(h)}$ be the oracle quantities defined in (1) Subsection 3.1 and $d_i^*$ the output of the MC Nonlinear Shrinkage algorithm with input*

---

**Algorithm 3.1** MC Nonlinear Shrinkage

---

1: Inputs: $\tilde{\sigma}, \tilde{\lambda}_1, \cdots, \tilde{\lambda}_n$ and a positive integer $K$.
2: **for** $k = 1, \cdots, K$ **do**
3:     Generate $\hat{Z}_k \sim \tilde{\sigma} n^{-1/2} GOE(n)$.
4:     Find the eigenvectors $\hat{g}_{1,k}, \cdots, \hat{g}_{n,k}$ of $diag(\tilde{\lambda}_1, \cdots, \tilde{\lambda}_n) + \hat{Z}_k$ such that $\hat{g}_{i,k}$ corresponds to the $i$–th largest eigenvalue.
5:     Set $\hat{d}_{i,k} = \hat{g}_{i,k}^\top diag(h(\tilde{\lambda}_1), \cdots, h(\tilde{\lambda}_n)) \hat{g}_{i,k}$.
6: Output: $d_i^* = K^{-1} \sum_{k=1}^K \hat{d}_{i,k}, 1 \le i \le n$.

---

$\tilde{\sigma}, \tilde{\lambda}_1, \cdots, \tilde{\lambda}_n, K \ge 1$. *Assume that* $\tilde{\sigma} \to \sigma$ *and* $n^{-1} \sum_{i=1}^n \delta_{\tilde{\lambda}_i} \xrightarrow{a.s.} H$. *Then, for any* $a, b \in [0, 1]$:

$$\frac{\sum_{[na]}^{[nb]} d_i^{(h)}}{n} - \frac{\sum_{[na]}^{[nb]} d_i^*}{n} \xrightarrow{a.s.} 0.$$

3.3. **Different Loss Functions.** So far we have been interested in the case of Frobenius loss. For some applications other losses might be more suitable. For this reason we shortly present how our results can be used to derive the optimal nonlinear shrinkage for some other choices of losses. Some of the losses we consider here (and many others) were studied for spiked covariance models in Donoho et al. [2018]. Below we will be interested in the following losses:

(1) Stein loss: $L^{st}(A, B) = tr(A^{-1}B - I) - \log \det A^{-1}B$.
(2) Divergence Loss: $L^{div}(A, B) = tr(A^{-1}B - I) + tr(B^{-1}A - I)$.
(3) The loss $L(A, B) = \left\| A^{-1}B - I \right\|_F^2$.

**Proposition 3.1.** *Assume that (using the notation from the Assumptions in Section 1)* $h_1 > 0$. *For any positive and bounded continuous function* $f$ *defined on an open set that eventually contains the eigenvalues of* $\hat{A}$ *we have almost surely:*

*(1) For the Stein loss* $L^{st}(A, f(\hat{A}))$ *we have:*

$$\lim_{n \to \infty} n^{-1} L^{st}(A, f(\hat{A})) = \int f(x) \frac{v_{1/t}(x)}{\pi} dx + \int \log t\, dH(t) - \int \log f(x) \frac{v(x)}{\pi} dx.$$

*This is minimized for* $f(x) = v(x)/v_{1/t}(x) = 1/f_{1/t}^*(x)$.

*(2) For the Stein loss* $L^{st}(f(\hat{A}), A)$ *we have:*

$$\lim_{n \to \infty} n^{-1} L^{st}(f(\hat{A}), A) = \int \frac{1}{f(x)} \frac{v_t(x)}{\pi} dx + \int \log f(x) \frac{v(x)}{\pi} dx - \int \log t\, dH(t) - 1.$$

*This is minimized for* $f(x) = v_t(x)/v(x) = f_t^*(x)$.

*(3) For the divergence loss* $L^{div}(A, f(\hat{A}))$ *we have:*

$$\lim_{n \to \infty} n^{-1} L^{div}(A, f(\hat{A})) = \int f(x) \frac{v_{1/t}(x)}{\pi} dx + \int \frac{1}{f(x)} \frac{v_t(x)}{\pi} dx - 2.$$

*This is minimized for* $f(x) = \sqrt{v_t(x)/v_{1/t}(x)} = \sqrt{f_t^*(x)/f_{1/t}^*(x)}$.

*(4) For the loss* $L(A, f(\hat{A}))$ *we have:*

$$\lim_{n \to \infty} n^{-1} L(A, f(\hat{A})) = 1 - 2 \int f(x) \frac{v_{1/t}(x)}{\pi} dx + \int f^2(x) \frac{v_{1/t^2}(x)}{\pi} dx.$$

*This is minimized for* $f(x) = v_{1/t}(x)/v_{1/t^2}(x) = f_{1/t}^*(x)/f_{1/t^2}^*(x)$.

*(5) For the loss $L(f(\hat{A}), A)$ we have:*

$$\lim_{n\to\infty} n^{-1} L(f(\hat{A}), A) = 1 - 2 \int \frac{1}{f(x)} \frac{v_t(x)}{\pi} dx + \int \frac{1}{f^2(x)} \frac{v_{t^2}(x)}{\pi} dx.$$

*This is minimized for $f(x) = v_{t^2}(x)/v_t(x) = f_{t^2}^*(x)/f_t^*(x)$.*

## 4. RECOVERY OF THE LIMITING SPECTRAL DISTRIBUTION

So far we have assumed the we know $\sigma, H$. In practice this is rarely true. Here we explain how those can be consistently estimated. First of all, assume that $\sigma$ is known. If $\sigma$ is unknown, we are going to see shortly that the problem can be ill-posed and further assumptions are needed to guarantee recovery of $\sigma, H$.

4.1. **Spectrum Recovery: known noise level.** When $\sigma$ is known, we suggest the procedure in Algorithm 4.1 that uses an optimization problem for recovering the eigenvalues of $A$.

---
**Algorithm 4.1** Population Eigenvalues Recovery

---
1: Inputs: $\hat{\lambda}_1, \cdots, \hat{\lambda}_n, \sigma$.
2: Sample $\hat{Z} \sim \sigma GOE(n)$.
3: For $T = (t_1, \cdots, t_n)^\top \in \mathbb{R}^n$ with $t_1 \geq \cdots t_n$, denote by $\hat{t}_1 \geq \cdots \geq \hat{t_n}$ the eigenvalues of $\mathrm{diag}(T) + n^{-1/2}\hat{Z}$.
4: Solve the optimization problem $T^* = \mathrm{argmin}_T\, n^{-1} \sum_{j=1}^n (\hat{t}_j - \hat{\lambda}_j)^2$.
5: Output $T^*$.

---

To minimize the objective above we suggest using the BFGS algorithm. A reasonable choice of a starting point that suggest is a point with independent Gaussian coordinates centered at the sample mean of the spectral distribution of $\hat{A}$. The optimization can be done quickly due to the fact that the gradients of the loss are easy to find in closed form. In particular, we have the following immediate proposition, which shows that the spectral decomposition of $T + \hat{Z}$ contains all the essential information to perform a BFGS update:

**Proposition 4.1.** *Using the notation from Algorithm 4.1, if $T + \hat{Z} = \sum_{j=1}^n \hat{t}_j \hat{x}_j \hat{x}_j^\top$ is the spectral decomposition of $T + n^{-1/2}\hat{Z}$, we have for all $i = 1, \cdots, n$:*

$$\partial_{t_i} \hat{t}_j = \hat{x}_{ij}^2.$$

*By $\hat{x}_{ij}$ we denote the $i$-th coordinate of $\hat{x}_j \in \mathbb{R}^n$.*

*Proof.* Let $E_{ii} \in \mathbb{R}^{n\times n}$ be the diagonal matrix with $i$-entry 1 and all other entries 0. Let $M_i(s) = T + n^{-1/2}\hat{Z} + sE_{ii}$. We have $\frac{d}{ds}M(s) = E_{ii}$, so using the *Hadamard first variation formula* (Page 57, Tao [2012]), we get

$$\partial_{t_i}\hat{t}_j = \hat{x}_j^\top E_{ii}\hat{x}_j = \hat{x}_{ij}^2.$$

$\square$

We have the following results that justify using this procedure:

**Theorem 4.1.** *Under the assumptions from Section 1, we have:*

*(1)*

$$\min_T \frac{1}{n} \sum_{i=1}^n \left(\hat{t}_i - \hat{\lambda}_i\right)^2 \xrightarrow{a.s.} 0.$$

*(2) If $T^*$ is a minimizer of the optimization problem above with $t_1^* \geq \cdots \geq t_n^*$, then*

$$\frac{1}{n} \sum_{i=1}^n (t_i^* - \lambda_i)^2 \xrightarrow{a.s.} 0.$$

**Remark 4.1.**     *(1) In the optimization problem we use only one copy of $\hat{Z} \sim \sigma GOE(n)$. In the high-dimensional limit $n \to \infty$ this is enough. Alternatively, as a regularization step, we could use multiple copies and solve the optimization problem repeatedly, getting solutions $T_1^*, \cdots, T_K^*$. We can then return $T^* = K^{-1} \sum_{i=1}^K T_i^*$.*

*(2) Theorem 4.1 shows that $n^{-1} \sum_{i=1}^n \delta_{t_i^*} \xrightarrow{a.s.} H$. This implies that the estimated eigenvalues can be used as input to Algorithm 3.1 and the assumptions of Theorem 3.2 will be satisfied.*

4.2. **Spectrum Recovery: unknown noise level.** If $\sigma$ is unknown, it is impossible to recover the measure $H$ simply by observing the free additive convolution with a semicircular measure of variance $\sigma^2$. To see why, assume that $H$ is semicircular with variance $s^2$. Then, $\mu_{H,\sigma^2}$ is semicircular with variance $s^2 + \sigma^2$ and it is impossible to separate the semicircular components of this measure. We conclude that further assumptions are needed. In fact, it is clear from the discussion above that only probability measures that cannot be written as the free additive convolution of a semicircular distribution and another probability measure are candidates for exact asymptotic recovery. For this reason, we are going to impose the following assumption throughout this section.

**Assumption 4.1.** *The measure $H$ cannot be written as the free additive convolution of a semicircular distribution with positive variance and a probability measure.*

In that case, if we solve the optimization problem from 4.1 for a choice $\hat{\sigma} < \sigma$, Theorem 4.1 suggests that the output will recover $H \boxplus \rho_{sc;\sigma^2 - \hat{\sigma}^2}$, while the objective should converge to 0. If we solve for a choice $\hat{\sigma} > \sigma$, then it is impossible to make the objective tend to 0. In particular, we have the following:

**Proposition 4.2.** *Let $R_n(\hat{\sigma})$ be the optimal value of the objective of the optimization problem in Algorithm 4.1 with $\sigma$ substituted by $\hat{\sigma}$. Then:*

*(1) For $\hat{\sigma} < \sigma$, $\limsup R_n(\hat{\sigma}) = 0$.*
*(2) For $\hat{\sigma} > \sigma$, $\liminf R_n(\hat{\sigma}) > 0$.*

Proposition 4.2 indicates that we can use a scree plot -type method to determine the noise level $\sigma$. In particular, we can solve the problem for several choices of the noise level and choose $\sigma$ before the objective becomes significantly larger than 0. This is going to be illustrated in Section 6.

## 5. ASYMPTOTIC EXPANSIONS

We study the asymptotic expansions of the oracle quantities and the optimal shrinkage functions in the regimes of "large noise" ($\sigma \to \infty$) and "small noise" ($\sigma \to 0$).

5.1. **The Large Noise Asymptotics.** If $\sigma \to \infty$, we have the following:

**Proposition 5.1.** *(1) If $Z \sim GOE(n)$, the oracle quantities $d_i^{(h)}$ defined in (1) almost surely satisfy:*

$$\lim_{n\to\infty} \lim_{\sigma\to\infty} \max_{1\le i\le n} \left| d_i^{(h)} - \int h(t)dH(t) \right| = 0.$$

*(2) The optimal shrinkage $f_h^*(x)$ satisfies*

$$\lim_{\sigma\to\infty} f_h^*(\sigma x) = \int h(t)dH(t)$$

*for $|x| < 2$.*

**Remark 5.1.** *Proposition 5.1 shows that in the regime of very large $\sigma$, the optimal nonlinear shrinkage quantities for estimation of $h(A)$ in Frobenius norm are essentially constant and achieve mean-squared-error equal to $\text{Var}\left[h(H)\right]$. This is reasonable, as an extremely large $\sigma$ should make estimation of $h(A)$ extremely hard. Notice that for $\sigma \to \infty$ the eigenvalues of $\hat{A}$ scale almost linearly with $\sigma$ and the limiting spectral distribution of $\sigma^{-1}\hat{A}$ is the semicircle law, which is indeed supported on $\left[-2, 2\right]$.*

5.2. **The Small Noise Asymptotics.** We now study the regime $\sigma \to 0$. Since for $\sigma = 0$ the eigenvectors of $A$ may not be uniquely determined, we assume for simplicity in this subsection that $A$ has distinct eigenvalues. In that case we have for $\sigma \to 0$:

**Proposition 5.2.** *If $h \in C^1(\mathbb{R})$ and $Z \sim GOE(n)$:*

*(1) The oracle quantities for $\sigma \to 0$ satisfy:*

$$\lim_{\sigma\to 0} \max_{1\le i\le n} \frac{\left| d_i^{(h)} - h(\lambda_i) \right|}{\sigma} = 0.$$

*(2)*

$$\lim_{n\to\infty} \lim_{\sigma\to 0} \frac{\left\| h(\hat{A}) - h(A) \right\|_F^2}{n\sigma^2} = \iint \frac{(h(t) - h(s))^2}{(t-s)^2} dH(t)dH(s).$$

We see from Proposition 5.2 that the mean-squared-error grows sublinearly in $\sigma^2$ for the optimal nonlinear shrinkage, if $\sigma$ is small, while using no shrinkage gives mean squared error $\approx \sigma^2 \iint (h(t) - h(s))^2/(t-s)^2 dH(t)dH(s)$ for $\sigma$ small. This is because, as we see from part 1 of Proposition 5.2 in the Gaussian case, the oracle quantities converge to $h(\lambda_i)$ fast for $\sigma \to 0$.

## 6. Numerical Experiments

6.1. **Experiments for Algorithm 4.1.** Here we consider three examples.

**Example 1.** Firstly we check the effectiveness of the deconvolution algorithm (Algorithm 4.1). For $H = (\delta_1 + \delta_4 + \delta_9)/3, \sigma^2 = 1$ and 20 equally spaced values of $n$ (starting from $n = 50$ and ending with $n = 1000$) we solve the optimization problem described in Algorithm 4.1. We start from 10 randomly initialized points and keep the stationary point of the objective that leads to the smallest value. We plot in
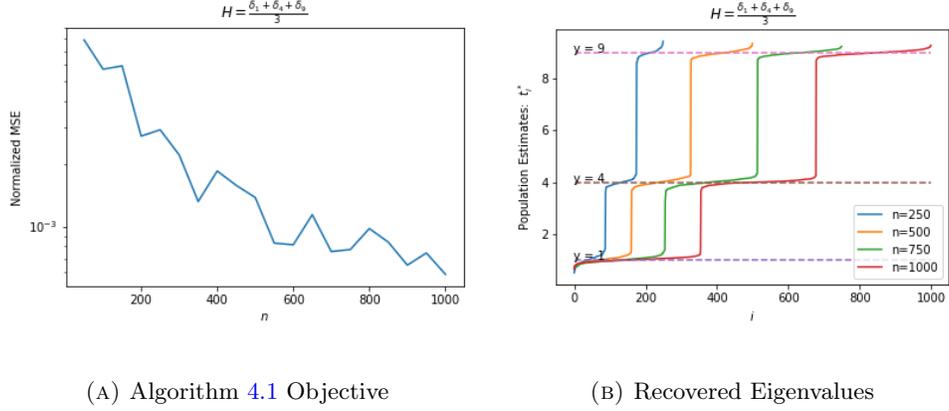
(A) Algorithm 4.1 Objective                    (B) Recovered Eigenvalues

FIGURE 1. Algorithm 4.1 Experiment for $H = (\delta_1 + \delta_4 + \delta_9)/3$

Figure 1, as a function of $n$, the resulting normalized mean squared error, which we define as

$$\frac{\frac{1}{n}\sum_{i=1}^{n}(t_i^* - \lambda_i)^2}{\mathrm{Var}[H]}.$$

We also present the recovered eigenvalues $t_i^*$ versus $i = 1, \cdots, n$ for the values $n = 250, 500, 750, 1000$.

**Example 2.** For a more complicated choice of spectral distribution $H$ we design the following experiment. We consider 200 randomly sampled points from circles centered at 0 with radii 0.5 and 1 respectively (presented with red and blue dots in the plot below). We add Gaussian noise with standard deviation 0.05 to the data. After generating those points, labeled as $x_1, \cdots, x_{200} \in \mathbb{R}^2$, we build the connectivity matrix $A \in \mathbb{R}^{200 \times 200}$ using the Gaussian kernel:

$$A_{ij} = \exp\left(-\frac{\|x_i - x_j\|^2}{2h^2}\right).$$

Here we choose $h = 0.1$. We assume that we have access only to a matrix

$$\hat{A} = A + \sqrt{\frac{2}{200}}Z,$$

where $Z$ is a standard Gaussian Wigner matrix. This corresponds to the choice $\sigma^2 = 2$. We use Algorithm 4.1 to estimate the eigenvalues of $A$. Below we plot the sample eigenvalues (that is the eigenvalues of $\hat{A}$), the true eigenvalues of $A$ and, finally, the estimated eigenvalues from the deconvolution algorithm. We see in Figure 2 that the reconstruction is very close.

**Example 3.** Finally, we consider an example with unknown noise level $\sigma^2$. In particular, we consider $H = (\delta_5 + \delta_{10})/2, \sigma^2 = 1, n = 200$. We take $Z$ to have entries drawn from a Laplace distribution. This time $\sigma^2$ is unknown, so we have to use several choices $\hat{\sigma}$ in the optimization problem and choose the largest $\hat{\sigma}$ for which the objective is close to 0. Figure 3 indicates using $\hat{\sigma}^2$ from 0.951 to 1.029. Refining the grid can give us an even closer estimate. We solve the optimization problem for $\hat{\sigma}^2 = 0.99$, which is the midpoint between the two values of $\sigma^2$ from above.
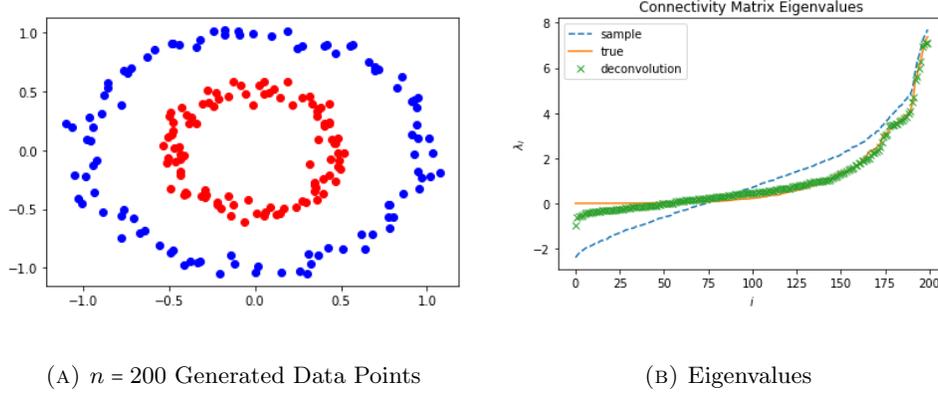
(A) $n = 200$ Generated Data Points

(B) Eigenvalues

FIGURE 2. Algorithm 4.1 Experiment for the Connectivity Matrix
created using a Gaussian Kernel with $h = 0.1$.



(B) Eigenvalues for the estimated $\hat{\sigma}^2 = 0.99$
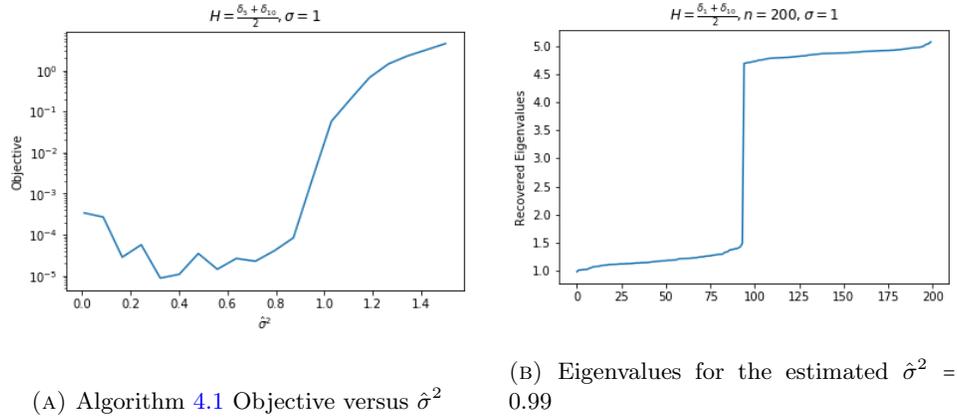
(A) Algorithm 4.1 Objective versus $\hat{\sigma}^2$

FIGURE 3. Algorithm 4.1 Experiment for $H = (\delta_1 + \delta_{10})/2, \sigma = 1$.
Here $\sigma$ is unknown and is estimated from $\hat{A}$.

6.2. **Noisy Linear Systems of Equations.** The first application we consider is
the following. We want to solve a linear system of equations of the form $Ax = b$,
whose solution we denote $x^* = A^{-1}b$. The matrix $A$ is unknown. Instead we have
access to a noisy estimate $\hat{A} = A + \sigma n^{-1/2} Z$, where $A, Z$ satisfy the assumptions
from Section 1. Solving $\hat{A}x = b$ gives $x = \hat{A}^{-1}b$. The problem is that $\hat{A}^{-1}$ might be
a very bad estimate of $A^{-1}$ and ill-conditioned. For this reason we suggest using
$x^{(f)} = f(\hat{A})b$, where $f$ is a bounded continuous function on $[h_1, h_2]$. Our goal is to
choose $f$ to minimize

$$\lim_{n \to \infty} \frac{\left\| x^{(f)} - x^* \right\|^2}{n}.$$

We study two different distributional assumptions on $b$.

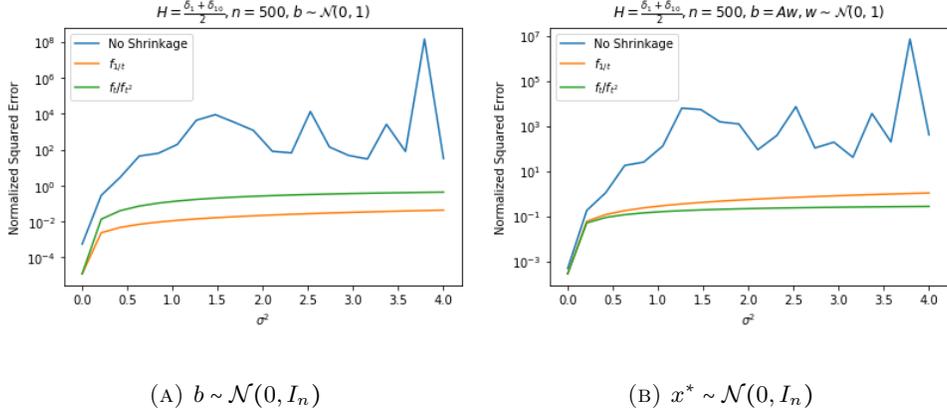(1) $b \sim \mathcal{N}(0, I_n)$. In that case using Lemma 7.2 we see that:

(A) $b \sim \mathcal{N}(0, I_n)$                    (B) $x^* \sim \mathcal{N}(0, I_n)$

FIGURE 4. $\left\| x^{(f)} - x^* \right\|^2 / \left\| x^* \right\|^2$ for different choices of $f$ and $H = (\delta_1 + \delta_{10})/2$.

$$\lim_{n \to \infty} \frac{\left\| x^{(f)} - x^* \right\|^2}{n} = \lim_{n \to \infty} \frac{\left\| (f(\hat{A}) - A^{-1})b \right\|^2}{n} = \lim_{n \to \infty} \frac{\left\| f(\hat{A}) - A^{-1} \right\|_F^2}{n}, \tag{8}$$

which is minimized for $f(x) = f_{1/t}^*(x)$.

(2) $b = Ax^*, x^* \sim \mathcal{N}(0, I_n)$. Similarly using Lemma 7.2 we see that:

$$\lim_{n \to \infty} \frac{\left\| x^{(f)} - x^* \right\|^2}{n} = \lim_{n \to \infty} \frac{\left\| (f(\hat{A})A - I_n)x^* \right\|^2}{n} = \lim_{n \to \infty} \frac{\left\| f(\hat{A})A - I_n \right\|_F^2}{n}. \tag{9}$$

Using exactly the same argument as in Proposition 3.1 we see that the limit is almost surely

$$1 - 2 \int f(x) \frac{v_t(x)}{\pi} dx + \int f^2(x) \frac{v_{t^2}(x)}{\pi} dx,$$

which is minimized for $f(x) = f_t^*(x)/f_{t^2}^*(x)$.

For $H = (\delta_1 + \delta_{10})/2, n = 500$ we plot in Figure 4 the normalized mean-squared-error (which we define as $\left\| x^{(f)} - x^* \right\|^2 / \left\| x^* \right\|^2$) for several values of $\sigma^2$.

We repeat the experiment for $H = (\delta_1 + \delta_4 + \delta_9)/3, n = 200$. The results can be seen in Figure 5.

In both cases we see that $\hat{A}$ becomes eventually ill-conditioned, if $\sigma^2$ increases. As expected, if $b \sim \mathcal{N}(0, I_n)$ the first shrinkage outperforms the second at all noise levels, while for $x^* \sim \mathcal{N}(0, I_n)$ the opposite is true.

### 6.3. Experiments for Algorithm 3.1.

We consider the problem of estimating $A, A^{-1}$ and $\sqrt{A}$ in Frobenius norm, when we only have access to $\hat{A}$. For $n = 500$ and several values of $\sigma$ we generate $\hat{A} = A + \sigma n^{-1/2}Z$, where $Z$ is a standard Gaussian Wigner matrix. Here $A$ is chosen as a diagonal matrix with diagonal entries chosen uniformly at random from $\{1, 4, 9\}$. We plot for $h(t) = t, h(t) = 1/t, h(t) = \sqrt{t}$ the oracle error and the error that can be achieved by using Algorithm 4.1 to recover the eigenvalues of $A$ and Algorithm 3.1 with $K = 1$ to perform nonlinear shrinkage. We see that in all cases the error achieved by our algorithm is very close to the
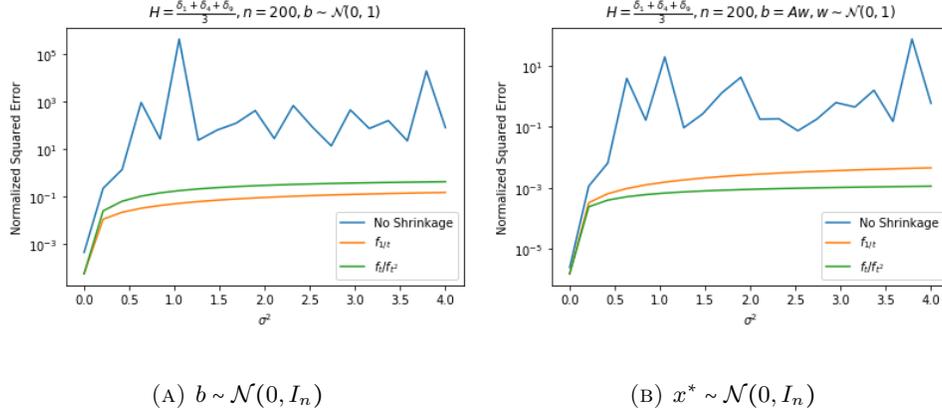
(A) $b \sim \mathcal{N}(0, I_n)$

(B) $x^* \sim \mathcal{N}(0, I_n)$

FIGURE 5. $\left\| x^{(f)} - x^* \right\|^2 / \left\| x^* \right\|^2$ for different choices of $f$ and $H = (\delta_1 + \delta_4 + \delta_9)/3$.



(A) Algorithm 3.1 for $h(t) = t$.

(B) Algorithm 3.1 for $h(t) = 1/t$.

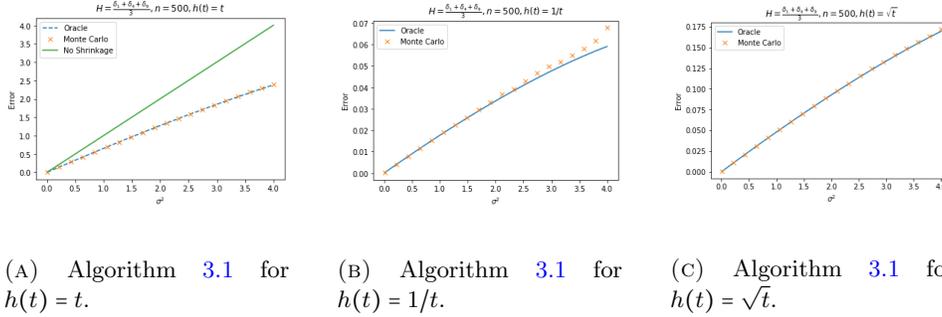(C) Algorithm 3.1 for $h(t) = \sqrt{t}$.

FIGURE 6. Algorithm 3.1 Experiment for $H = (\delta_1 + \delta_4 + \delta_9)/3$. We plot the oracle and shrinkage errors versus $\sigma$ for $h(t) = t, 1/t, \sqrt{t}$.

oracle. For $\sigma^2$ large, notice that the problem of eigenvalue recovery for $A$ becomes increasingly harder, hence the error in estimation of $\lambda_1, \cdots, \lambda_n$ increases. This can lead to problems for the function $h(t) = 1/t$ which is unbounded near 0, hence we clip all the recovered eigenvalues that we get from Algorithm 4.1 to be at least 0.3. Notice that for $h(t) = 1/t, h(t) = \sqrt{t}$ we do not plot the no shrinkage

## 7. PROOFS

7.1. **Proofs for Section 2.** We start by presenting a well-known lemma for the tails of a standard Gaussian random variable.

**Lemma 7.1.** *For any $M > 0$ and $Z \sim \mathcal{N}(0,1)$ we have*

$$\mathbb{P}(|Z| > M) \le 2M^{-1} \exp\left(-x^2/2\right).$$

*Proof.* We have

$$
\mathbb{P}(|Z| > M) = 2\mathbb{P}(Z > M) = 2 \int_M^\infty \exp\left(-x^2/2\right) dx \tag{10}
$$

$$
\leq 2 \int_M^\infty \frac{x}{M} \exp\left(-x^2/2\right) dx = 2M^{-1} \exp\left(-x^2/2\right).
$$

$\square$

We will need the following lemma which is adapted from Lemma 7.8, Lemma 7.9 and Lemma 7.10 from Erdős and Yau [2017].

**Lemma 7.2.** *Let $q \geq 2$ and $X_1, \cdots, X_N, Y_1, \cdots, Y_N$ be independent random variables with mean 0, variance 1 and $2q$-th moment bounded by $c_0$. Then, for any deterministic $(b_i)_{1 \leq i \leq N}, (a_{ij})_{1 \leq i,j \leq N}$ we have for some positive constant $C_q = C_q(c_0)$:*

$$
\left\| \sum_i b_i(X_i^2 - 1) \right\|_q \leq C_q \left( \sum_i |b_i|^2 \right)^{\frac{1}{2}} \tag{11}
$$

$$
\left\| \sum_{i,j} a_{ij} X_i Y_j \right\|_q \leq C_q \left( \sum_{i,j} a_{ij}^2 \right)^{\frac{1}{2}} \tag{12}
$$

$$
\left\| \sum_{i \neq j} a_{ij} X_i X_j \right\|_q \leq C_q \left( \sum_{i \neq j} a_{ij}^2 \right)^{\frac{1}{2}} \tag{13}
$$

*Proof of Theorem 2.2.* The proof involves two main steps.
  (1) **Step 1:** Show that the theorem holds for if $Z \sim GOE(n)$.
  (2) **Step 2:** Reduce the problem to the case of bounded random variables as entries of $Z$.
  (3) **Step 3:** Show that the results are universally true and independent of the distribution of $Z$ as long as $Z, A$ are asymptotically free.

**Step 1:** For $Z \sim GOE(n)$ (which is invariant under conjugation by an orthogonal matrix) it is enough to consider the case of diagonal matrix $A$. If $A = \mathrm{diag}(A_1, \cdots, A_n)$, then using the Schur complement formula we have that the $i$-th diagonal entry of $h(A) \left( \hat{A} - z \right)^{-1}$ is given by

$$
h(A_i) / \left( A_i - z - \sigma n^{-1/2} Z_{ii} - \sigma^2 n^{-1} a_i^\top \left( A_{-i} + n^{-1/2} Z_{-i} - z \right)^{-1} a_i \right).
$$

Here $A_{-i}$ is the $(n-1) \times (n-1)$ matrix that we get if we omit the $i$-th element $A_i$ of $A$, $a_i \in \mathbb{R}^{(n-1)}$ the $i$-th row of $Z$ if we omit the diagonal element.

We have from Lemma 7.1 for any fixed $\epsilon > 0$ and $n \geq \epsilon^{-2}$:

$$
\mathbb{P}\left( \max_{1 \leq i \leq n} |Z_{ii}| > \epsilon\sqrt{n} \right) \leq n\mathbb{P}\left( |Z_{ii}| > \epsilon\sqrt{n} \right) = \mathcal{O}\left( n \exp\left(-n\epsilon^2/2\right) \right).
$$

Since $\epsilon$ was arbitrary, we conclude by the Borel-Cantelli lemma that

$$
\max_{1 \leq i \leq n} |Z_{ii}| / \sqrt{n} \xrightarrow{a.s.} 0.
$$

Similarly, by Lemma 7.2 for $q = 3$,

$$
\max_{1 \leq i \leq n} n^{-1} \left| a_i^\top \left( A_{-i} + \sigma n^{-1/2} Z_{-i} - z \right)^{-1} a_i - tr\left( \left( A_{-i} + \sigma n^{-1/2} Z_{-i} - z \right)^{-1} \right) \right|.
$$

Using the Cauchy interlacing formula (Tao [2012]), we see that it must also be true that

$$\max_{1 \le i \le n} \left| n^{-1} tr \left( \left( A_{-i} + \sigma n^{-1/2} Z_{-i} - z \right)^{-1} \right) - m_{\hat{\mu}}(z) \right| \xrightarrow{a.s.} 0, \tag{14}$$

as for fixed $z \in \mathbb{C}^+$ the differences

$$\left| tr \left( \left( A_{-i} + \sigma n^{-1/2} Z_{-i} - z \right)^{-1} \right) - tr \left( \left( A + \sigma n^{-1/2} Z - z \right)^{-1} \right)(z) \right|$$

are going to be uniformly bounded (due to the interlacing phenomenon).

As a consequence, we see that

$$n^{-1} tr \left( h(A) \left( \hat{A} - z \right)^{-1} \right) = o(1) + \sum_{i=1}^{n} \frac{h(A_i)}{t_i - z - \sigma^2 m_{\hat{\mu}}(z)},$$

which proves the result for $z \in \mathbb{C}^+$ and $Z \sim GOE(n)$.

**Step 2:** Fix $M > 0$. For this step we assume, in order to slightly simplify the formulas, that without loss of generality that $\sigma = 1$. Define $Z_{i,j}^{(M)} = Z_{ij} \mathbb{I}_{|Z_{ij}| < M}$. We also define $\hat{A}^{(M)} = A + n^{-1/2} Z^{(M)}$. We have

$$\left| n^{-1} tr \left( h(A)(\hat{A}^{(M)} - z)^{-1} \right) - n^{-1} tr \left( h(A)(\hat{A} - z)^{-1} \right) \right|$$

$$= \frac{\left| tr \left( h(A)(\hat{A} - z)^{-1}(Z - Z^{(M)})(\hat{A}^{(M)} - z)^{-1} \right) \right|}{n \sqrt{n}}$$

$$= \frac{\left| tr \left( (\hat{A}^{(M)} - z)^{-1} h(A)(\hat{A} - z)^{-1}(Z - Z^{(M)}) \right) \right|}{n \sqrt{n}} \tag{15}$$

$$\le \frac{\left\| (\hat{A}^{(M)} - z)^{-1} h(A)(\hat{A} - z)^{-1} \right\|_{op} \left\| (Z - Z^{(M)}) \right\|_F}{n}.$$

Notice that here we have used the fact that for two $n \times n$ matrices $M_1, M_2$ we have

$$|tr (M_1 M_2)| \le \|M_1\|_F \|M_2\|_F \le \sqrt{n} \|M_1\|_{op} \|M_2\|_F ,$$

where the first inequality follows from Cauchy-Schwartz in $\mathbb{R}^{n \times n}$ and the second one from the fact that the Frobenius norm of a real matrix is the $l_2$−norm of its singular values.

We conclude from (15) that for a fixed bounded continuous function and a fixed complex number $z$ in the upper half-plane we have

$$\left| n^{-1} tr \left( h(A)(\hat{A}^{(M)} - z)^{-1} \right) - n^{-1} tr \left( h(A)(\hat{A} - z)^{-1} \right) \right| = \mathcal{O} \left( n^{-1} \left\| Z - Z^{(M)} \right\|_F \right). \tag{16}$$

We now observe that

$$\mathbb{E} \left[ \left\| (Z - Z^{(M)}) \right\|_F^2 \right] = n^2 \mathbb{E} \left[ Z_{11}^2; |Z_{11}| \ge M \right]$$

and

$$\mathrm{Var} \left[ \left\| Z - Z^{(M)} \right\|_F^2 \right] = \mathcal{O} \left( n^2 \mathbb{E} \left[ Z_{11}^4; |Z_{11}| > M \right] \right).$$

Fix any $\epsilon > 0$ and take $M$ large enough such that $\mathbb{E}[Z_{11}^2; |Z_{11}| > M] \le \epsilon^2/2$ and $\mathbb{E} \left[ Z_{11}^4; |Z_{11}| > M \right] \le 1$. Then we have

$$\mathbb{P} \left( n^{-1} \left\| Z - Z^{(M)} \right\|_F > \epsilon \right) = \mathbb{P} \left( n^{-2} \left\| Z - Z^{(M)} \right\|_F^2 > \epsilon^2 \right)$$

$$\le \mathbb{P} \left( n^{-2} \left\| Z - Z^{(M)} \right\|_F^2 - \mathbb{E} \left[ n^{-2} \left\| Z - Z^{(M)} \right\|_F^2 \right] > \epsilon^2/2 \right) = \mathcal{O} \left( n^{-2} \epsilon^{-4} \right). \tag{17}$$

Using the Borel-Cantelli lemma we conclude that. almost surely, $n^{-1} \left\| Z - Z^{(M)} \right\|_F \leq \epsilon$ eventually. Using (15) we see that for $M$ large enough we have eventually almost surely

$$
\begin{aligned}
&\left| n^{-1} tr \left( h(A)(\hat{A}^{(M)} - z)^{-1} \right) - n^{-1} tr \left( h(A)(\hat{A} - z)^{-1} \right) \right| \\
&\leq \epsilon \left\| (\hat{A}^{(M)} - z)^{-1} h(A)(\hat{A} - z)^{-1} \right\|_{op} \leq \epsilon Im(z)^{-2} \left\| h \right\|_\infty.
\end{aligned}
\tag{18}
$$

To finish this step, we define $\mu_M = \mathbb{E}\left[ Z_{ij}^{(M)} \right], \sigma_M = \sqrt{\mathrm{Var}\left[ Z_{ij}^{(M)} \right]}$ and $\tilde{Z}_{ij} = \left( Z_{ij}^{(M)} - \mu_M \right)/\sigma_M$, which are random variables with mean 0 and variance 1. A similar argument shows that, if $\tilde{A} = A + n^{-1/2}\tilde{Z}$ and $M$ is large enough, then eventually almost surely we have

$$
\left| n^{-1} tr \left( h(A) \left( \hat{A}^{(M)} - z \right)^{-1} \right) - n^{-1} tr \left( h(A) \left( \tilde{A} - z \right)^{-1} \right) \right| \leq \epsilon.
\tag{19}
$$

To see why, using the same bound as in (15), we see that

$$
\begin{aligned}
&\left| n^{-1} tr \left( h(A) \left( \hat{A}^{(M)} - z \right)^{-1} \right) - n^{-1} tr \left( h(A) \left( \tilde{A} - z \right)^{-1} \right) \right| \leq \\
&\frac{\left\| \left( \hat{A}^{(M)} - z \right)^{-1} h(A) \left( \tilde{A} - z \right)^{-1} \right\|_{op} \left\| \tilde{Z} - Z^{(M)} \right\|_F}{n} \leq n^{-1} \left\| h \right\|_\infty Im(z)^{-2} \left\| \tilde{Z} - Z^{(M)} \right\|_F.
\end{aligned}
\tag{20}
$$

It remains to bound $\tilde{Z} - Z^{(M)}$ in Frobenius norm. Let $e = (1, \cdots, 1)^\top \in \mathbb{R}^n$. Then,

$$
\tilde{Z} - Z^{(M)} = \tilde{Z}^{(M)}(1 - \sigma_M^{-1}) - \frac{\mu_M}{\sigma_M} ee^\top.
$$

Using this we get

$$
\begin{aligned}
\left\| \tilde{Z} - Z^{(M)} \right\|_F^2 &= (1 - \sigma_M^{-1})^2 \left\| Z^{(M)} \right\|_F^2 - 2\frac{\mu_M}{\sigma_M}(1 - \sigma_M^{-1}) e^\top Z^{(M)} e + n^2 \frac{\mu_M^2}{\sigma_M^2} \\
&\leq n^2 \frac{\mu_M^2}{\sigma_M^2} - 2\frac{\mu_M}{\sigma_M}(1 - \sigma_M^{-1}) e^\top Z^{(M)} e + (1 - \sigma_M^{-1}) \left( \left\| Z \right\|_F + \left\| Z - Z^{(M)} \right\|_F \right)^2
\end{aligned}
\tag{21}
$$

Now we know that:
(1) $\lim_{M \to \infty} \mu_M = 0$ and $\lim_{M \to \infty} \sigma_M = 1$ from the dominated convergence theorem.
(2) $n^{-1} \left\| Z - Z^{(M)} \right\|$ can be made arbitrarily small eventually (by choosing $M$ large enough), using (17) and the Borel-Cantelli lemma.
(3) $n^{-1} \left\| Z \right\|_F \leq \sqrt{n^{-1} \left\| Z \right\|_{op}}$, which converges to 2 almost surely (Tao [2012]).
(4) Finally,

$$
n^{-2} e^\top Z^{(M)} e = \frac{\sum_{1 \leq i,j \leq n} Z_{ij} \mathbb{I}_{|Z_{ij}| < M}}{n^2} \xrightarrow{a.s.} \mathbb{E}\left[ Z_{11}; |Z_{11}| < M \right]
$$

by the strong law of large numbers.

Taking all of the above into consideration, we see that for $M$ large enough we see that for $M$ large enough we have eventually almost surely that the bound from (19) is true. Since $\tilde{Z}$ is a Wigner ensemble with bounded entries, we have reduced the problem to the case of bounded random variables.

**Step 3:** We now show that under the assumptions in Subsection 1.1 the theorem is also true. The idea is to show, using free-probabilistic tools, that if $\tilde{Z}$ is

a Wigner ensemble with all moments finite, the the limit of the trace functionals of interest depends on the noise distribution via only its first two moments. First of all, notice that it is enough to prove the result for a polynomial $h$ and the extend to a general continuous function by a simple density argument. As a result, it is enough to consider $h(t) = t^k, k \in \mathbb{N}$ and show that $n^{-1} tr \left( A^k \left( \hat{A} - z \right)^{-1} \right)$ has a limit almost surely and the limit does not depend on the distribution of $Z$. Similarly, it is enough to show that for any $m \in \mathbb{N}$ the trace functional $n^{-1} tr \left( A^k \hat{A}^m \right)$ has a limit almost surely and the limit does not depend on the distribution of $Z$. Writing $\hat{A}^m = \left( A + \sigma n^{-1/2} Z \right)^m$ and expanding in monomial terms we see that $n^{-1} tr \left( A^k \hat{A}^m \right)$ is the sum of a finite number of terms all of which have the form $n^{-1} tr \left( A^{n_1} (\sigma n^{-1/2} Z)^{m_1} \cdots A^{n_s} \left( n^{-1/2} Z \right)^{m_s} \right)$ for some $s \geq 1$ and nonnegative integers $n_1, m_1, \cdots, n_s, m_s$. Since $Z$ is a Wigner ensemble and $A$ is independent of $Z$, we conclude that $A, n^{-1/2} Z$ are almost surely asymptotically free (Theorem 20 in Mingo and Speicher [2017]). As a consequence, we have that all terms of the form $n^{-1} tr \left( A^{n_1} (\sigma n^{-1/2} Z)^{m_1} \cdots A^{n_s} \left( n^{-1/2} Z \right)^{m_s} \right)$ converge almost surely and the limit depends only on the limiting spectral distributions of $A, n^{-1/2} Z$, which are given by $H$ and a semicircular distribution respectively. In particular, the limit is independent of the distribution of $Z$. We conclude that the limit is the same as with the Gaussian assumption on $Z$. This completes the proof. $\qquad \square$

### 7.2. Proofs for Section 3.

*Proof of Theorem 3.1.* First of all, let $f$ be an analytic function on the complex plane. Then, we have using Cauchy's integral formula:

$$\frac{tr \left( f(\hat{A}) h(A) \right)}{n} = -\frac{1}{2\pi i} \oint_{|z|=R} \frac{tr \left( h(A) \left( \hat{A} - z \right)^{-1} \right)}{n} f(z) dz,$$

where the integral is considered on a fixed circle centered at $0$ with radius $R$ such that eventually $\left\| \hat{A} \right\|_{op} \leq R/2$. Consider

$$R_n(z) = \frac{tr \left( h(A) \left( \hat{A} - z \right)^{-1} \right)}{n}$$

and for a fixed $z$ let us denote by $R_\infty(z)$ the almost sure limit of $R_n(z)$ described in Theorem 2.2. Then, for $n$ large enough we have almost surely that:

$$|R_n(z)| = \left| \frac{tr \left( h(A) \left( \hat{A} - z \right)^{-1} \right)}{n} \right| \leq \frac{\|h(A)\|_{op}}{(R - R/2)} = 2R^{-1} \|h(A)\|_{op},$$

and

$$|R_n'(z)| = \left| \frac{tr \left( h(A) \left( \hat{A} - z \right)^{-2} \right)}{n} \right| \leq \frac{\|h(A)\|_{op}}{(R - R/2)^2} = 4R^{-2} \|h(A)\|_{op},$$

so on the circle $\{z \in \mathbb{C} : |z| = R\}$ the sequence of functions $\{R_n(z)\}_{n \geq 1}$ almost surely consists of functions that are uniformly bounded and equicontinuous. Fix some

$\epsilon > 0$ and consider a finite subset $C \subset \{z \in \mathbb{C} : |z| = R\}$ such that for any $z \in \mathbb{C}$ with $|z| = R$ there exists $\tilde{z} \in C$ such that $|z - \tilde{z}| < \epsilon$. Since for any such $z, \tilde{z}$ we have

$$|R_n(z) - R_\infty(z)| \le |R_n(z) - R_n(\tilde{z})| + |R_n(\tilde{z}) - R_\infty(\tilde{z})| + |R_\infty(\tilde{z}) - R_\infty(z)|$$

$$= \mathcal{O}\left(\epsilon + \sup_{\tilde{z} \in C} |R_n(\tilde{z}) - R_\infty(\tilde{z})|\right),$$

we know that almost surely

$$\limsup_{|z|=R} \sup |R_n(z) - R_\infty(z)| = \mathcal{O}(\epsilon).$$

Since $\epsilon$ was arbitrary we conclude that $R_n \xrightarrow{a.s.} R_\infty$ uniformly on $\{z \in \mathbb{C} : |z| = R\}$. Using this result we see that

$$\frac{tr\left(h(A)f\left(\hat{A}\right)\right)}{n} \xrightarrow{a.s.} -\frac{1}{2\pi i} \oint_{|z|=R} R_\infty(z)f(z)dz = -\frac{1}{2\pi i} \oint_{\Gamma_\delta} R_\infty(z)f(z)dz,$$

where we have changed the integral to be on a counterclockwise curve $\Gamma_\delta$ which we take to be a rectangle with vertices $\pm R \pm i\delta$. Taking $\delta \downarrow 0$ we get

$$-\frac{1}{2\pi i} \oint R_\infty(z)f(z)dz$$
$$= -\frac{1}{2\pi i} \lim_{\delta \downarrow 0} \left[\int_{-R}^{R} R_\infty(x - i\delta)f(x - i\delta)dx - \int_{-R}^{R} R_\infty(x + i\delta)f(x + i\delta)\right] \qquad (22)$$
$$= \int \frac{v_h(x)}{\pi} f(x)dx.$$

In other words, we have shown that

$$\frac{tr\left(f(\hat{A})h(A)\right)}{n} \xrightarrow{a.s.} \int \frac{v_h(x)}{\pi} f(x)dx$$

for $f$ analytic. Using a simple density argument we see that this result is actually true for any function $f$ that is continuous and bounded in an open set that contains the support of $\mu_{H,\sigma^2}$ and eventually all the eigenvalues of $\hat{A}$.

We now see that

$$\frac{\left\|f(\hat{A}) - h(A)\right\|_F^2}{n} = \frac{\left\|f(\hat{A})\right\|_F^2 - 2tr\left(f(\hat{A}h(A)) + \|h(A)\|_F^2\right)}{n} \qquad (23)$$
$$\xrightarrow{a.s.} \int f^2(x)\frac{v(x)}{\pi}dx - 2\int \frac{v_h(x)}{\pi}f(x)dx + \int h^2(x)dH(x).$$

Minimizing over $f$ we see that for $x$ in the support of $\mu_{H,\sigma^2}$ the minimizer satisfies $f^*(x) = v_h(x)/v(x)$. This completes the proof. $\qquad \square$

*Proof of Corollary 3.1.*     (1) We have

$$\int \frac{tdH(t)}{t - z - \sigma^2 m_{\hat{\mu}}(z)} = 1 + \int \frac{z + \sigma^2 m_{\hat{\mu}}(z)}{t - z - \sigma^2 m_{\hat{\mu}}(z)} = 1 + (z + \sigma^2 m_{\hat{\mu}}(z))m_{\hat{\mu}}(z),$$

where the last equality follows from Proposition 2.1.

(2) Using

$$\frac{1}{t(t-z-\sigma^2 m_{\hat{\mu}}(z))} = \frac{1}{z+\sigma^2 m_{\hat{\mu}}(z)}\left[\frac{1}{t-z\sigma^2 m_{\hat{\mu}}(z)} - \frac{1}{t}\right],$$

we have

$$\int \frac{dH(t)}{t(t-z-\sigma^2 m_{\hat{\mu}}(z))} = \frac{m_{\hat{\mu}}(z)-m_H(0)}{z+\sigma^2 m_{\hat{\mu}}(z)},$$

so the formula for the asymptotically optimal shrinkage for $A^{-1}$ follows.

(3) We have

$$
\begin{aligned}
&\int \frac{t^2 dH(t)}{t-z-\sigma^2 m_{\hat{\mu}}(z)}\\
&= \int \frac{t^2 - (z+\sigma^2 m_{\hat{\mu}}(z))^2}{t-z-\sigma^2 m_{\hat{\mu}}(z)} dH(t) + (z+\sigma^2 m_{\hat{\mu}}(z))^2 \int \frac{dH(t)}{t-z-\sigma^2 m_{\hat{\mu}}(z)}\\
&= \int \left(t+z+\sigma^2 m_{\hat{\mu}}(z)\right) dH(t) + (z+\sigma^2 m_{\hat{\mu}}(z))^2 m_{\hat{\mu}}(z)\\
&= \int t dH(t) + z + \sigma^2 m_{\hat{\mu}}(z) + (z+\sigma^2 m_{\hat{\mu}}(z))^2 m_{\hat{\mu}}(z).
\end{aligned}
\tag{24}
$$

$\square$

*Proof of Theorem 3.2.* First of all, we observe that it is enough to prove the theorem for $K = 1$, where $\hat{Z}_1 \sim \tilde{\sigma} n^{-1/2} GOE(n)$. Hence we consider only that case and ignore the dependency on $k$ in the subscripts in Algorithm 3.1. We will write $\tilde{\Lambda} = \mathrm{diag}(\tilde{\lambda}_1, \cdots, \tilde{\lambda}_n)$ for the diagonal matrix in Step 2 of Algorithm 3.1. We will denote by $\tilde{m}_1 \geq \cdots \geq \tilde{m}_n$ the eigenvalues of $\tilde{\Lambda} + \hat{Z}_1$. From Theorem 2.2 we know that:

$$n^{-1}\sum_{i=1}^n d_i^* \delta_{\tilde{m}_i} \xrightarrow{a.s.} \mu_h,$$

where $\mu_h$ is a finite measure with Stieltjes transform given by

$$\int \frac{h(t) dH(t)}{t-z-\sigma^2 m_{\hat{\mu}}(z)}.$$

In addition, we know that $n^{-1}\sum_{i=1}^n \delta_{\tilde{m}_i}$ converges weakly almost surely to the additive free convolution of $H$ with a semicircular distribution with variance $\sigma^2$, which is a probability measure $H \boxplus \rho_{sc;\sigma^2}$ without atoms (Biane [1997]). We conclude that for any $x_1, x_2 \in \mathbb{R}$ we have $n^{-1}\sum_{i=1}^n d_i^* \mathbb{I}_{\tilde{m}_i \in [x_1,x_2]} \xrightarrow{a.s.} \mu_h([x_1,x_2])$. Similarly $n^{-1}\sum_{i=1}^n d_i^{(h)} \mathbb{I}_{\hat{\lambda}_i \in [x_1,x_2]} \xrightarrow{a.s.} \mu_h([x_1,x_2])$. Since $H \boxplus \rho_{sc;\sigma^2}$ has no atoms, the proof is completed if we consider $x_1, x_2$ be the $a, b$-quantiles respectively of $\mu_{H,\sigma^2}$.

$\square$

*Proof of Proposition 3.1.* In the proof of Theorem 3.1 we saw that, if $f$ satisfies the assumptions of Proposition 3.1, then almost surely

$$\lim_{n\to\infty} n^{-1} tr\left(f(\hat{A})h(A)\right) = \int f(x)\frac{v_h(x)}{\pi} dx. \tag{25}$$

(1)

$$n^{-1}L^{st}(A, f(\hat{A})) = \frac{1}{n}tr\left(A^{-1}f(\hat{A})\right) - 1 + \frac{1}{n}\sum_{i=1}^{n}\log\lambda_i - \frac{1}{n}\sum_{i=1}^{n}\log f(\hat{\lambda}_i)$$

$$\xrightarrow{a.s.} \int f(x)\frac{v_{1/t}(x)}{\pi}dx - 1 + \int \log t\,dH(t) - \int \log f(x)\frac{v(x)}{\pi}dx,$$

(26)

where we used (25) and the fact that the spectrum of $A$ converges weakly almost surely to $H$, while the spectrum of $\hat{A}$ converges weakly almost surely to the measure $H \boxplus \rho_{sc;\sigma^2}$ with density $v(x)/\pi$. Minimizing the integrand with respect to $f(x)$ for $x$ fixed is straightforward using derivatives and gives the desired result.

(2)

$$n^{-1}L^{st}(f(\hat{A}), A) = \frac{1}{n}tr\left(f(\hat{A})^{-1}A\right) - 1 + \frac{1}{n}\sum_{i=1}^{n}\log f(\hat{\lambda}_i) - \frac{1}{n}\sum_{i=1}^{n}\log\lambda_i$$

$$\xrightarrow{a.s.} \int \frac{1}{f(x)}\frac{v_t(x)}{\pi}dx - \int \log t\,dH(t) + \int \log f(x)\frac{v(x)}{\pi}dx - 1.$$

(27)

Minimizing with respect to $f$ is again straightforward.

(3) Using (25) we get:

$$n^{-1}L^{div}(A, f(\hat{A})) = \frac{tr\left(Af(\hat{A})^{-1}\right)}{n} + \frac{tr\left(A^{-1}f(\hat{A})\right)}{n} - 2$$

$$\xrightarrow{a.s.} \int f(x)\frac{v_{1/t}(x)}{\pi}dx + \int \frac{1}{f(x)}\frac{v_t(x)}{\pi}dx - 2.$$

(28)

(4)

$$n^{-1}\left\|A^{-1}f(\hat{A}) - I\right\|_F^2 = \frac{tr\left(A^{-2}f(\hat{A})^2\right)}{n} - 2\frac{tr\left(A^{-1}f(\hat{A})\right)}{n} + 1$$

$$\xrightarrow{a.s.} \int f^2(x)\frac{v_{1/t^2}(x)}{\pi}dx - 2\int f(x)\frac{v_{1/t}(x)}{\pi}dx + 1.$$

(29)

(5)

$$n^{-1}\left\|Af(\hat{A})^{-1} - I\right\|_F^2 = \frac{tr\left(A^2f(\hat{A})^{-2}\right)}{n} - 2\frac{tr\left(Af(\hat{A})^{-1}\right)}{n} + 1$$

$$\xrightarrow{a.s.} 1 - 2\int \frac{1}{f(x)}\frac{v_t}{\pi}dx + \int \frac{1}{f^2(x)}\frac{v_{t^2}(x)}{\pi}dx.$$

(30)

$\square$

### 7.3. Proofs for Section 4.

*Proof of Theorem 4.1.*     (1) We take $t_i = \lambda_i$. Then,

$$n^{-1}\sum_{i=1}^{n}\delta_{\hat{t}_i} \xrightarrow{a.s.} H \boxplus \rho_{sc;\sigma^2}.$$

In addition,

$$n^{-1}\sum_{i=1}^{n}\delta_{\hat{\lambda}_i} \xrightarrow{a.s.} H \boxplus \rho_{sc;\sigma^2}.$$

Finally, applying Weyl's inequality ((1.54) in Tao [2012]) to $\hat{A} = A + \sigma n^{-1/2}Z$ and using the fact that for any $\epsilon > 0$ the eigenvalues of $n^{-1/2}Z$

eventually lie in $[-2 - \epsilon, 2 + \epsilon]$ almost surely, we have $\lambda_i - 2\sigma + o(1) \leq \hat{\lambda}_i \leq \lambda_i + 2\sigma + o(1)$, so that $\hat{\lambda}_i$ are almost surely uniformly bounded. Similarly for $\hat{t}_i$. We conclude that for this choice of $t_i$'s the 2-Wasserstein distance of

$$n^{-1} \sum_{i=1}^{n} (\hat{t}_i - \hat{\lambda}_i)^2 \xrightarrow{a.s.} 0.$$

The proof is completed.

(2) If we denote by $\nu_n = n^{-1} \sum_{i=1}^{n} \delta_{t_i^*}$ the probability measure that corresponds to the solution to the optimization problem in Algorithm 4.1, then we know that $\nu_n$ is tight sequence of probability measures. To see why, by Weyl's eigenvalue inequality and the fact that almost surely the largest eigenvalue of $n^{-1/2}Z$ tends to 2 and the largest eigenvalue of $n^{-1/2}Z$ tends to -2 we get $t_i^* - 2\sigma + o(1) \leq \hat{t}_i^* \leq t_i^* + 2\sigma + o(1)$, so

$$\left| t_i^* - \hat{\lambda}_i \right| \leq \left| \hat{t}_i^* - \hat{\lambda}_i \right| + 2\sigma + o(1). \tag{31}$$

For any large $M > 0$ fixed we now see from (31) that

$$M^2 \frac{\left| \{ i : 1 \leq i \leq n, |t_i - \lambda_i| > M \} \right|}{n} \leq n^{-1} \sum_{i=1}^{n} (t_i - \hat{\lambda}_i)^2 \mathbb{I}\{ |t_i - \hat{\lambda}_i| > M \}$$

$$\leq n^{-1} \left( \sqrt{\sum_{i=1}^{n} (\hat{t}_i^* - \hat{\lambda}_i)^2} + 2\sigma\sqrt{n} + o(\sqrt{n}) \right)^2 = (R_n(\sigma) + 2\sigma + o(1))^2. \tag{32}$$

On the other hand, we know that any weak subsequential limit of $\nu_n$ has to be equal to $H$, as from the previous part of the theorem we know that $n^{-1} \sum_{i=1}^{n} \delta_{t_i^*} \boxplus \rho_{sc;\sigma^2}$ converges weakly to $H \boxplus \rho_{sc;\sigma^2}$. We conclude that $\nu_n \xrightarrow{\mathcal{D}} H$ almost surely. Fix $\epsilon > 0$ and consider $M > 0$ large enough (to be determined later). For the moment we assume that $[-M/2, M/2]$ contains $[h_1, h_2]$ from Assumption 3 in section 1. In addition, we assume that $|\hat{\lambda}_i| \leq M/2$ for all $i$.

Then, using the triangle inequality we have:

$$\sqrt{n^{-1} \sum_{i=1}^{n} (t_i^*)^2 \mathbb{I}\{ |t_i^*| > M \}} \leq \sqrt{n^{-1} \sum_{i=1}^{n} (\hat{t}_i^* - \hat{\lambda}_i)^2} + \sqrt{n^{-1} \sum_{i=1}^{n} \hat{\lambda}_i^2 \mathbb{I}\{ |t_i^*| > M \}}$$

$$+ \sqrt{n^{-1} \sum_{i=1}^{n} (\hat{t}_i^* - t_i^*)^2 \mathbb{I}\{ |t_i^*| > M \}}$$

$$\leq \sqrt{n^{-1} \sum_{i=1}^{n} (\hat{t}_i^* - \hat{\lambda}_i)^2} + (\max_{1 \leq i \leq n} |\hat{\lambda}_i| + 2\sigma + o(1)) \sqrt{\frac{\left| i : 1 \leq i \leq n, |t_i^*| > M \right|}{n}}$$

$$\leq \sqrt{n^{-1} \sum_{i=1}^{n} (\hat{t}_i^* - \hat{\lambda}_i)^2} + (\max_{1 \leq i \leq n} |\hat{\lambda}_i| + 2\sigma + o(1)) \sqrt{\frac{\left| i : 1 \leq i \leq n, |t_i^* - \hat{\lambda}_i| > M/2 \right|}{n}} \tag{33}$$

The first term in the above inequality goes to 0 almost surely, as we saw in part (1), while the second term can be made arbitrarily small from the bound in (31). So if we choose $M$ large enough, then eventually almost surely we have

$$n^{-1} \sum_{i=1}^{n} (t_i^*)^2 \mathbb{I}\{ |t_i^*| > M \} \leq \epsilon. \tag{34}$$

Combining the bound in (34), since $\epsilon$ was arbitrary, with the fact that $\nu_n \xrightarrow{\mathcal{D}} H$ and $n^{-1} \sum_{i=1}^n \delta_{\lambda_i} \xrightarrow{\mathcal{D}} H$ we see that

$$n^{-1} \sum_{i=1}^n (t_i^* - \lambda_i)^2 \xrightarrow{a.s.} 0.$$

$\square$

*Proof of Proposition 4.2.* (1) If $\hat{\sigma} < \sigma$, then taking $t_i$ to be the $(i-1)/p$ quantile of $H \boxplus \rho_{sc;\sigma^2 - \hat{\sigma}^2}$ (which is the additive free-convolution of $H$ with a semicircular distribution with variance $\sigma^2 - \hat{\sigma}^2$) gives that the empirical distribution $n^{-1} \sum_{i=1}^n \hat{t}_i$ converges weakly almost surely to

$$H \boxplus \rho_{sc;\sigma^2 - \hat{\sigma}^2} \boxplus \rho_{sc;\hat{\sigma}^2} = H \boxplus \rho_{sc;\sigma^2}.$$

As a consequence, the Wasserstein 2-distance of $n^{-1} \sum_{i=1}^n \hat{t}_i$ and $n^{-1} \sum_{i=1}^n \hat{\lambda}_i$ converges almost surely to 0, so

$$n^{-1} \sum_{i=1}^n (\hat{t}_i - \hat{\lambda}_i)^2 \xrightarrow{a.s.} 0.$$

This shows that almost surely $\limsup R_n(\hat{\sigma}) = 0$.

(2) Fix $\hat{\sigma} > \sigma$ and consider the event $\mathcal{A} = \{\liminf R_{\hat{\sigma}} = 0\}$. Assume that $A$ has positive probability. Then, for an $\omega \in \mathcal{A}$ there exists a sequence $n_k \uparrow \infty$ such that $R_{n_k}(\hat{\sigma}; \omega) \to 0$.

If we denote by $\nu_n = n^{-1} \sum_{i=1}^n \delta_{t_i^*(\omega)}$ the probability measure that corresponds to the solution to the optimization problem in Algorithm 4.1, then we know that $\nu_{n_k}$ is tight sequence of probability measures, as in the proof of Theorem 4.1. We conclude that there exists a subsequence of $\{\nu_{n_k}\}_{k \geq 1}$ that converges weakly to a probability measure $\tilde{H}$. Then, we must have, due to the fact that $R_{n_k;\omega}(\hat{\sigma}) \to 0$,

$$\tilde{H} \boxplus \rho_{sc;\hat{\sigma}^2} = H \boxplus \rho_{sc;\sigma^2}.$$

Since $\tilde{H} \boxplus \rho_{sc;\hat{\sigma}^2} = \tilde{H} \boxplus \rho_{sc;\hat{\sigma}^2 \boxplus \sigma^2} \boxplus \rho_{sc;\sigma^2}$, we deduce that $H = \tilde{H} \boxplus \rho_{sc;\hat{\sigma}^2 - \sigma^2}$. This is a contradiction, so $\mathbb{P}(\mathcal{A}) = 0$ and the proof is completed. $\square$

## 7.4. **Proofs for Section 5.**

*Proof of Proposition 5.1.* (1) For $n$ fixed and $\sigma \to \infty$, $\sigma^{-1} \hat{A}_n = \sigma^{-1} A_n + n^{-1/2} Z_n \to n^{-1/2} Z_n$ and the eigenvectors of $\hat{A}_n$ converge to the eigenvectors of $Z_n$ which are uniformly distributed with respect to the Haar measure. Let us denote by $z_1, \cdots, z_n$ the $l_2$−normalized eigenvectors of $Z$. We have

$$\mathbb{E}\left[(z_i^\top w_j)^2 | A\right] = \frac{1}{n} \Rightarrow \mathbb{E}\left[z_i^\top A z_i | A\right] = \frac{tr(A)}{n}.$$

Applying Theorem 5.1.4 in Vershynin [2018] for the function $f(X) = X^\top A X$ (which is Lipschitz on the unit sphere with Lipschitz constant $2\|A\|_{op}$) we have that there exists a constant $C > 0$ such that for any $\epsilon > 0$ and any $i$:

$$\mathbb{P}\left(\left|z_i^\top A z_i - n^{-1} tr(A)\right| > \epsilon\right) \leq 2 \exp\left(-\frac{cn\epsilon^2}{\|A\|_{op}^2}\right) \leq 2 \exp\left(-\frac{cn\epsilon^2}{h_2^2}\right). \tag{35}$$

Using the union bound we have

$$\mathbb{P}\left(\max_{1\le i\le n}\left|z_i^\top A z_i - n^{-1}tr\,(A)\right| > \epsilon\right) \le 2n\exp\left(-\frac{cn\epsilon^2}{h_2^2}\right). \tag{36}$$

The Borel-Cantelli lemma implies that almost surely we eventually have

$$\max_{1\le i\le n}\left|z_i^\top A z_i - n^{-1}tr\,(A)\right| \le \epsilon.$$

Since $\epsilon$ was arbitrary we have

$$\max_{1\le i\le n}\left|z_i^\top A z_i - n^{-1}tr\,(A)\right| \xrightarrow{a.s} 0.$$

To finish the proof we see that

$$n^{-1}tr(A) = n^{-1}\sum_{1\le i\le n}h(\lambda_i) \xrightarrow{a.s.} \int h(t)dH(t).$$

(2) We have from Theorem 3.1 that for $|x| < 2$ (such that $\sigma x$ eventually lies in the support of $H \boxplus \rho_{sc;\sigma^2}$):

$$f_h^*(\sigma x) = \int \frac{h(t)dH(t)}{(\sigma^{-1}t - x - \sigma u(\sigma x))^2 + \sigma^2 v(\sigma x)^2}. \tag{37}$$

Firstly, we will show that $(x + \sigma u(\sigma x))^2 + \sigma^2 v(\sigma x)^2 \to 1$.

Let $\alpha(z;\sigma) = \sigma m(\sigma z), \tilde{\alpha}(x) = \lim_{\epsilon\downarrow 0}\alpha(x + i\epsilon)$. Then, we have

$$\tilde{\alpha}(x) = \int \frac{dH(t)}{\sigma^{-1}t - x - \tilde{\alpha}(x)},$$

so $\tilde{\alpha(x)}$ converges, as $\sigma \to \infty$, to the solution of the equation $\beta = -1/(x+\beta)$ that lies on the upper half plane. Notice that $\beta$ is the limit of the Stieltjes transform of the semicircular distribution with variance 1 on the real axis. So

$$\beta^2 + x\beta + 1 = 0 \to (\beta + x)^2 - (\beta + x)x + 1 = 0 \to |\beta + x| = 1.$$

This is exactly what we claimed, in particular that

$$(x + \sigma u(\sigma x))^2 + \sigma^2 v(\sigma x)^2 \xrightarrow{\sigma\to\infty} 1.$$

The rest will follow from Scheffé's lemma. In particular, we have for any $x \in \mathbb{R}$ from the equation that defines $m_{\hat\mu}$:

$$\int \frac{dH(t)}{(\sigma^{-1}t - x - \sigma u(\sigma x))^2 + \sigma^2 v(\sigma x)^2} = 1.$$

As a consequence, for each $x$ the measure

$$\frac{dH(t)}{(\sigma^{-1}t - x - \sigma u(\sigma x))^2 + \sigma^2 v(\sigma x)^2}$$

is a probability measure that converges weakly to $H$ as $\sigma \downarrow 0$. The proof is completed.

$\square$

*Proof of Proposition 5.2.* (1) We have

$$\hat{w}_i^\top h(A)\hat{w}_i = w_i^\top h(A)w_i + 2\sigma\frac{d\hat{w}_i}{d\sigma}\Big|_{\sigma=0}h(A)w_i + o(\sigma)$$

$$= w_i^\top h(A)w_i + 2h(\lambda_i)\sigma w_i^\top\frac{d\hat{w}_i}{d\sigma}\Big|_{\sigma=0} + o(\sigma). \tag{38}$$

Since $\hat{w}_i^\top \hat{w}_i = 1$, we get

$$w_i^\top \frac{d\hat{w}_i}{d\sigma}\big|_{\sigma=0} = 0.$$

We conclude that

$$\lim_{\sigma\to 0} \max_{1\le i\le n} \frac{\left|\hat{w}_i^\top h(A)\hat{w}_i - h(\lambda_i)\right|}{\sigma} = 0.$$

(2) From the Hadamard variation formulas for the eigenvalues and eigenvectors of $\hat{A}$ (Erdős and Yau [2017]) we know that:

$$\frac{d\hat{\lambda}_i}{d\sigma}\big|_{\sigma=0} = w_i^\top \frac{Z}{\sqrt{n}} w_i$$

$$\frac{d\hat{w}_i}{d\sigma}\big|_{\sigma=0} = \sum_{j\ne i} \frac{\frac{w_i^\top Z w_j}{\sqrt{n}}}{\lambda_i - \lambda_j} w_j$$

Using these we have for $n$ fixed and $\sigma \to 0$ we have under the convention $h'(x) = (h(x) - h(y))/(x - y)$ for $x = y$ :

$$h(\hat{A}) = \sum_{i=1}^n h(\hat{\lambda}_i)\hat{w}_i\hat{w}_i^\top = \sum_{i=1}^n h\left(\lambda_i + \sigma w_i^\top \frac{Z}{\sqrt{n}} w_i\right)\hat{w}_i\hat{w}_i^\top + o(\sigma)$$

$$= \sum_{i=1}^n \left[h(\lambda_i) + \sigma w_i^\top \frac{Z}{\sqrt{n}} w_i h'(\lambda_i)\right]\hat{w}_i\hat{w}_i^\top + o(\sigma)$$

$$= \sum_{i=1}^n \left[h(\lambda_i) + \sigma \frac{w_i^\top Z w_i}{\sqrt{n}} h'(\lambda_i)\right]\left[w_i w_i^\top + \sigma \sum_{j\ne i} \frac{\frac{w_i^\top Z w_j}{\sqrt{n}}}{\lambda_i - \lambda_j}(w_i w_j^\top + w_j w_i^\top)\right] + o(\sigma) \quad (39)$$

$$= \sum_{i=1}^n h(\lambda_i) w_i w_i^\top + \frac{\sigma}{2}\sum_{i,j=1}^n \frac{w_i^\top Z w_j}{\sqrt{n}}\frac{h(\lambda_i) - h(\lambda_j)}{\lambda_i - \lambda_j}(w_i w_j^\top + w_j w_i^\top) + o(\sigma)$$

$$= h(A) + \frac{\sigma}{2}\sum_{i,j=1}^n \frac{w_i^\top Z w_j}{\sqrt{n}}\frac{h(\lambda_i) - h(\lambda_j)}{\lambda_i - \lambda_j}(w_i w_j^\top + w_j w_i^\top) + o(\sigma).$$

This gives us

$$\lim_{\sigma\to 0}\frac{\left\|h(\hat{A}) - h(A)\right\|_F^2}{n\sigma^2} = \frac{1}{n}\sum_{i,j=1}^n \left(\frac{w_i^\top Z w_j}{\sqrt{n}}\right)^2 \frac{(h(\lambda_i) - h(\lambda_j))^2}{(\lambda_i - \lambda_j)}.$$

Due to the rotational invariance of $Z$ we can assume that $u_i$ is the $i$-th standard basis vector. We get

$$\lim_{\sigma\to 0}\frac{\left\|h(\hat{A}) - h(A)\right\|_F^2}{n\sigma^2} = n^{-2}\sum_{i,j=1}^n Z_{ij}^2\left(\frac{h(\lambda_i) - h(\lambda_j)}{\lambda_i - \lambda_j}\right)^2. \tag{40}$$

Writing

$$m_{ij} = \left(\frac{h(\lambda_i) - h(\lambda_j)}{\lambda_i - \lambda_j}\right)^2,$$

we know that $m_{ij} \le \|h'\|_\infty^2$ and

$$n^{-2}\sum_{i,j=1}^n m_{ij} \xrightarrow{a.s.} \iint \frac{(h(t) - h(s))^2}{(t-s)^2}dH(t)dH(s).$$

In addition,

$$\mathbb{E}\left[n^{-2}\sum_{i,j=1}^{n}Z_{ij}^2\frac{(h(\lambda_i)-h(\lambda_j)^2}{(\lambda_i-\lambda_j)^2}|m_{ij},1\le i,j\le n\right]=n^{-2}\sum_{i,j=1}^{n}m_{ij}^2+n^{-2}\sum_{i=1}^{n}m_{ii}^2,$$

so

$$\mathbb{E}\left[n^{-2}\sum_{i,j=1}^{n}Z_{ij}^2\frac{(h(\lambda_i)-h(\lambda_j))^2}{(\lambda_i-\lambda_j)^2}|m_{ij},1\le i,j\le n\right] \tag{41}$$

$$\xrightarrow{a.s.}\iint\frac{(h(t)-h(s))^2}{(t-s)^2}dH(t)dH(s)$$

Finally, we have

$$\mathrm{Var}\left[n^{-2}\sum_{i,j=1}^{n}m_{ij}Z_{ij}^2|m_{ij},1\le i,j\le n\right]=\mathcal{O}(n^{-2}),$$

so we get

$$n^{-2}\sum_{i,j=1}^{n}m_{ij}Z_{ij}^2-\mathbb{E}\left[n^{-2}\sum_{i,j=1}^{n}Z_{ij}^2\frac{(h(\lambda_i)-h(\lambda_j))^2}{(\lambda_i-\lambda_j)^2}|m_{ij},1\le i,j\le n\right]\xrightarrow{a.s.}0.$$

We deduce from (41) that

$$n^{-2}\sum_{i,j=1}^{n}m_{ij}Z_{ij}^2\xrightarrow{a.s.}\iint\frac{(h(t)-h(s))^2}{(t-s)^2}dH(t)dH(s).$$

$\square$

## References

James O Berger, Dongchu Sun, Chengyuan Song, et al. Bayesian analysis of the covariance matrix of a multivariate normal distribution with a new class of priors. *Annals of Statistics*, 48(4):2381–2403, 2020.

Philippe Biane. On the free convolution with a semi-circular distribution. *Indiana University Mathematics Journal*, pages 705–718, 1997.

Joël Bun, Romain Allez, Jean-Philippe Bouchaud, and Marc Potters. Rotational invariant estimator for general noisy matrices. *IEEE Transactions on Information Theory*, 62(12):7475–7490, 2016.

David L Donoho and Matan Gavish. The optimal hard threshold for singular values is $4/\sqrt{3}$, 2013.

David L Donoho, Matan Gavish, and Iain M Johnstone. Optimal shrinkage of eigenvalues in the spiked covariance model. *Annals of statistics*, 46(4):1742, 2018.

Noureddine El Karoui et al. Spectrum estimation for large dimensional covariance matrices using random matrix theory. *The Annals of Statistics*, 36(6):2757–2790, 2008.

László Erdős and Horng-Tzer Yau. *A dynamical approach to random matrix theory*, volume 28. American Mathematical Soc., 2017.

Philip A Etter and Lexing Ying. Operator augmentation for noisy elliptic systems. *arXiv preprint arXiv:2010.09656*, 2020.

William James and Charles Stein. Estimation with quadratic loss. In *Breakthroughs in statistics*, pages 443–460. Springer, 1992.

Weihao Kong, Gregory Valiant, et al. Spectrum estimation from samples. *Annals of Statistics*, 45(5):2218–2247, 2017.

Olivier Ledoit and Sandrine Péché. Eigenvectors of some large sample covariance matrix ensembles. *Probability Theory and Related Fields*, 151(1):233–264, 2011.

Olivier Ledoit and Michael Wolf. A well-conditioned estimator for large-dimensional covariance matrices. *Journal of multivariate analysis*, 88(2):365–411, 2004.

Olivier Ledoit and Michael Wolf. Spectrum estimation: A unified framework for covariance matrix estimation and pca in large dimensions. *Journal of Multivariate Analysis*, 139:360–384, 2015.

Olivier Ledoit, Michael Wolf, et al. Nonlinear shrinkage estimation of large-dimensional covariance matrices. *The Annals of Statistics*, 40(2):1024–1060, 2012.

Olivier Ledoit, Michael Wolf, et al. Analytical nonlinear shrinkage of large-dimensional covariance matrices. *Annals of Statistics*, 48(5):3043–3065, 2020.

Panagiotis Lolas. Regularization in high-dimensional regression and classification via random matrix theory. *arXiv preprint arXiv:2003.13723*, 2020.

James A Mingo and Roland Speicher. *Free probability and random matrices*, volume 35. Springer, 2017.

Sheehan Olver and Raj Rao Nadakuditi. Numerical computation of convolutions in free probability theory. *arXiv preprint arXiv:1203.1958*, 2012.

Marc Potters and Jean-Philippe Bouchaud. *A First Course in Random Matrix Theory: For Physicists, Engineers and Data Scientists*. Cambridge University Press, 2020.

N Raj Rao and Alan Edelman. The polynomial method for random matrices. *Foundations of Computational Mathematics*, 8(6):649–702, 2008.

Terence Tao. *Topics in random matrix theory*, volume 132. American Mathematical Soc., 2012.

Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.

Eugene P. Wigner. On the distribution of the roots of certain symmetric matrices. *Annals of Mathematics*, 67(2):325–327, 1958. ISSN 0003486X. URL http://www.jstor.org/stable/1970008.

Ruoyong Yang and James O Berger. Estimation of a covariance matrix using the reference prior. *The Annals of Statistics*, pages 1195–1211, 1994.