



# DIGITAL ACCESS TO SCHOLARSHIP AT HARVARD

## Bulk universality for generalized Wigner matrices

The Harvard community has made this article openly available.  
[Please share](#) how this access benefits you. Your story matters.

<b>Citation</b>	Erdos, László, Horng-Tzer Yau, and Jun Yin. 2011. "Bulk Universality for Generalized Wigner Matrices." <i>Probab. Theory Relat. Fields</i> 154 (1-2) (October 6): 341–407. doi:10.1007/s00440-011-0390-3. <a href="http://dx.doi.org/10.1007/s00440-011-0390-3">http://dx.doi.org/10.1007/s00440-011-0390-3</a> .
<b>Published Version</b>	<a href="https://doi.org/10.1007/s00440-011-0390-3">doi:10.1007/s00440-011-0390-3</a>
<b>Accessed</b>	March 18, 2018 2:14:19 AM EDT
<b>Citable Link</b>	<a href="http://nrs.harvard.edu/urn-3:HUL.InstRepos:25427234">http://nrs.harvard.edu/urn-3:HUL.InstRepos:25427234</a>
<b>Terms of Use</b>	This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Open Access Policy Articles, as set forth at <a href="http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#OAP">http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#OAP</a>

*(Article begins on next page)*

# Bulk universality for generalized Wigner matrices

László Erdős<sup>1\*</sup>, Horng-Tzer Yau<sup>2†</sup> and Jun Yin<sup>2 ‡</sup>

Institute of Mathematics, University of Munich,  
Theresienstr. 39, D-80333 Munich, Germany  
lerdos@math.lmu.de <sup>1</sup>

Department of Mathematics, Harvard University  
Cambridge MA 02138, USA  
htyau@math.harvard.edu, jyin@math.harvard.edu <sup>2</sup>

Aug 14, 2010

## Abstract

Consider  $N \times N$  Hermitian or symmetric random matrices  $H$  where the distribution of the  $(i, j)$  matrix element is given by a probability measure  $\nu_{ij}$  with a subexponential decay. Let  $\sigma_{ij}^2$  be the variance for the probability measure  $\nu_{ij}$  with the normalization property that  $\sum_i \sigma_{ij}^2 = 1$  for all  $j$ . Under essentially the only condition that  $c \leq N\sigma_{ij}^2 \leq c^{-1}$  for some constant  $c > 0$ , we prove that, in the limit  $N \rightarrow \infty$ , the eigenvalue spacing statistics of  $H$  in the bulk of the spectrum coincide with those of the Gaussian unitary or orthogonal ensemble (GUE or GOE). We also show that for band matrices with bandwidth  $M$  the local semicircle law holds to the energy scale  $M^{-1}$ .

**AMS Subject Classification (2010):** 15B52, 82B44

*Keywords:* Random band matrix, Local semicircle law, sine kernel.

---

\*Partially supported by SFB-TR 12 Grant of the German Research Council

†Partially supported by NSF grants DMS-0602038, 0757425, 0804279

‡Partially supported by NSF grants DMS-100165

# 1 Introduction

One key universal quantity for random matrices is the eigenvalue gap distribution. Although the density of eigenvalues may depend on the specific model, the gap distribution or the short distance correlation function are believed to depend only on the symmetry class of the ensembles but are otherwise independent of the details of the distributions. There are two types of universality: the edge universality and the bulk universality. In this paper, we will focus on the bulk universality concerning the interior of the spectrum. The bulk universality was proved for very general classes of invariant ensembles (see, e.g. [3, 6, 7, 8, 9, 25, 26, 27] and references therein). For non-invariant ensembles, in particular for matrices with i.i.d. entries (Wigner matrices), the bulk universality was difficult to establish due to the lack of an explicit expression for the joint distribution of the eigenvalues.

The first rigorous partial result for bulk universality in the non-unitary case was given by Johansson [23] (see also Ben Arous and Pécché [2] and the recent improvement [24]) stating that the bulk universality holds for Gaussian divisible *Hermitian* ensembles, i.e., Hermitian ensembles of the form

$$\widehat{H} + sV, \tag{1.1}$$

where  $\widehat{H}$  is a Wigner matrix,  $V$  is an independent standard GUE matrix and  $s$  is a positive constant of order one. The restriction on Gaussian divisibility turned out to be very difficult to remove. In a series of papers [12, 13, 14, 17], we developed a new approach to prove the universality. The first step was to derive the local semicircle law, an estimate of the local eigenvalue density, down to energy scales containing around  $\log N$  eigenvalues. Once such a strong form of the local semicircle law was obtained, the result of [23, 2] can be extended to a Gaussian convolution with variance only  $s^2 \asymp N^{-1+\varepsilon}$ . This tiny Gaussian component can then be removed via a reverse heat flow argument and this proves [17] the bulk universality for Hermitian ensembles provided that the distributions of the matrix elements are sufficiently differentiable.

The bulk universality for Hermitian ensembles was also proved later on by Tao and Vu [31] under the condition that the first four moments of the matrix elements match those of GUE, but without the differentiability assumption. The condition on the fourth moment was already removed in [31] by using the result for Gaussian divisible ensembles of [23, 2]; the third moment condition was then removed in [18] by using the result of [17].

The four moment theorem [31] is also valid for the symmetric ensembles, but the restriction on the matching of the first four moments cannot be weakened for the following reason. The key input to remove the fourth moment matching condition for the Hermitian case, the universality of the Gaussian divisible ensembles [23, 2], relied entirely on the asymptotic analysis of an *explicit* formula, closely related to a formula in Brézin-Hikami [5, 23], for the correlation functions of the eigenvalues for the *Hermitian ensembles*  $\widehat{H} + sV$ . Since similar formulas for symmetric matrices are very complicated, the corresponding result is not available and thus the matching of the fourth moment cannot be removed in this way. Although there is a proof [16] of universality for  $s^2 \geq N^{-3/4}$  without using this formula, the main ingredient of that proof, establishing the uniqueness of the local equilibria of the Dyson Brownian motion, still heavily used explicit formulas related to GUE.

In [15] a completely different strategy was introduced based on a *local relaxation flow*, which locally behaves like a Dyson Brownian motion, but has a faster decay to equilibrium. This approach entirely eliminates explicit formulas and it gives a unified proof for the universality of *symmetric* and Hermitian Wigner matrices [15]. It was further generalized [19] to quaternion self-dual Wigner matrices and sample covariance matrices. The method not only applies to all these specific ensembles, but it also gives a conceptual interpretation that the occurrence of the universality is due to the relaxation to local equilibrium of the DBM. We remark that very recently the results of [31] were also extended to sample covariance matrices [33].

The main input of all these methods [17, 15, 19] and [31, 33] is an estimate of the local density of eigenvalues, the local semicircle law. This has been developed in the previous work on Wigner matrices [12, 13, 14], where the matrix elements were i.i.d. random variables. In this paper, we extend this method to random matrices with independent, but not necessarily identically distributed entries. If we denote the variance of the  $(i, j)$  entry of the matrix by  $\sigma_{ij}^2$ , our main interest is the case that  $\sigma_{ij}$  are not a constant but they satisfy the normalization condition  $\sum_i \sigma_{ij}^2 = 1$  for all  $j$ . We will call such matrix ensembles *universal Wigner matrices*. For these ensembles Guionnet [21] and Anderson-Zeitouni [1] proved that the density of the eigenvalues converges to the Wigner semi-circle law. The simplest case is that of *generalized Wigner matrices*, where  $N\sigma_{ij}^2$  is uniformly bounded from above and below by two fixed positive numbers. In this case, we prove the local semicircle law down to essentially the smallest possible energy scale  $N^{-1}$  (modulo  $\log N$  factors). A much more difficult case is the *Wigner band matrices* where, roughly speaking,  $\sigma_{ij}^2 = 0$  if  $|i - j| > M$  for some  $M < N$ . In this case, we obtain the local semicircle law to the energy scale  $M^{-1}$ . We note that a certain three-dimensional version of Gaussian band matrices was considered by Disertori, Pinson and Spencer [10] using the supersymmetric method. They proved that the expectation of the density of eigenvalues is smooth and it coincides with the Wigner semicircle law.

With the local semicircle law proved up to the almost optimal scale, applying the method of [15, 19] leads to the identification of the correlation functions and the gap distribution for generalized Wigner matrices provided that the distribution of the matrix elements is continuous and satisfies the logarithmic Sobolev inequality. These additional assumptions can be removed if one can extend the Tao-Vu theorem [31] to generalized Wigner matrices. In Section 8, we will introduce an approach based on a Green's function comparison theorem, which states that the joint distributions of Green's functions of two ensembles at different energies with imaginary parts of order  $1/N$  are identical provided that the first three moments of the two ensembles coincide and the fourth moments are close. Since local correlation functions and the gap distribution of the eigenvalues can be identified from Green's functions, it follows that the local correlation functions of these two ensembles are identical at the scale  $1/N$ . We can thus use this theorem to remove all continuity and logarithmic Sobolev inequality restrictions in our approach. In particular, this leads to the bulk universality for generalized Wigner matrices with the subexponential decay being essentially the only assumption on the probability law. We note that one major technical difficulty in [31], the level repulsion estimate, is not needed in the proof of the Green's function comparison theorem. It will be clear in Section 8 that, once the local semicircle law is established, the Green's function comparison theorem is a simple consequence of the standard resolvent perturbation theory.

## 2 Main results

We now state the main results of this paper. Since all our results hold for both Hermitian and symmetric ensembles, we will state the results for Hermitian matrices only. The modifications to the symmetric case are straightforward and they will be omitted. Let  $H = (h_{ij})_{i,j=1}^N$  be an  $N \times N$  Hermitian matrix where the matrix elements  $h_{ij} = \overline{h_{ji}}$ ,  $i \leq j$ , are independent random variables given by a probability measure  $\nu_{ij}$  with mean zero and variance  $\sigma_{ij}^2$ . The variance of  $h_{ij}$  for  $i > j$  is  $\sigma_{ij}^2 = \mathbb{E}|h_{ij}|^2 = \sigma_{ji}^2$ . For simplicity of the presentation, we assume that for any fixed  $1 \leq i < j \leq N$ ,  $\text{Re } h_{ij}$  and  $\text{Im } h_{ij}$  are i.i.d. with distribution  $\omega_{ij}$  i.e.,  $\nu_{ij} = \omega_{ij} \otimes \omega_{ij}$  in the sense that  $\nu_{ij}(dh) = \omega_{ij}(d\text{Re } h)\omega_{ij}(d\text{Im } h)$ , but this assumption is not essential for the result. The distribution  $\nu_{ij}$  and its variance  $\sigma_{ij}^2$  may depend on  $N$ , but we suppress this in the notation. We assume that, for any  $j$  fixed,

$$\sum_i \sigma_{ij}^2 = 1. \tag{2.1}$$

Matrices with independent, zero mean entries and with the normalization condition (2.1) will be called *universal Wigner matrices*. For a forthcoming review on this matrix class, see [29], where the terminology of *random band matrices* was used.

Define  $C_{inf}$  and  $C_{sup}$  by

$$C_{inf} := \inf_{N,i,j} \{N\sigma_{ij}^2\} \leq \sup_{N,i,j} \{N\sigma_{ij}^2\} =: C_{sup}. \quad (2.2)$$

Note that  $C_{inf} = C_{sup}$  corresponds to the standard Wigner matrices and the condition  $0 < C_{inf} \leq C_{sup} < \infty$  defines more general Wigner matrices with comparable variances.

We will also consider an even more general case when  $\sigma_{ij}$  for different  $(i, j)$  indices are not comparable. The basic parameter of such matrices is the quantity

$$M := \frac{1}{\max_{i,j} \sigma_{ij}^2}. \quad (2.3)$$

A special case is the band matrix, where  $\sigma_{ij} = 0$  for  $|i - j| > W$  with some parameter  $W$ . In this case,  $M$  and  $W$  are related by  $M \leq CW$ .

Denote by  $B := \{\sigma_{ij}^2\}_{i,j=1}^N$  the matrix of variances which is symmetric and doubly stochastic by (2.1), in particular it satisfies  $-1 \leq B \leq 1$ . Let the spectrum of  $B$  be supported in

$$\text{Spec}(B) \subset [-1 + \delta_-, 1 - \delta_+] \cup \{1\} \quad (2.4)$$

with some nonnegative constants  $\delta_{\pm}$ . We will always have the following spectral assumption

$$1 \text{ is a simple eigenvalue of } B \text{ and } \delta_- \text{ is a positive constant, independent of } N. \quad (2.5)$$

The local semicircle law will be proven under this general condition, but the precision of the estimate near the spectral edge will also depend on  $\delta_+$  in an explicit way. For the orientation of the reader, we mention two special cases of universal Wigner matrices that provided the main motivation for our work.

*Example 1. Generalized Wigner matrix.* In this case we have

$$0 < C_{inf} \leq C_{sup} < \infty, \quad (2.6)$$

and one can easily prove that 1 is a simple eigenvalue of  $B$  and (2.4) holds with

$$\delta_{\pm} \geq C_{inf}, \quad (2.7)$$

i.e., both  $\delta_-$  and  $\delta_+$  are positive constants independent of  $N$ .

*Example 2. Band matrix.* The variances are given by

$$\sigma_{ij}^2 = W^{-1} f\left(\frac{[i-j]_N}{W}\right), \quad (2.8)$$

where  $W \geq 1$ ,  $f : \mathbb{R} \rightarrow \mathbb{R}_+$  is a bounded nonnegative symmetric function with  $\int f = 1$  and we defined  $[i-j]_N \in \mathbb{Z}$  by the property that  $[i-j]_N \equiv i-j \pmod{N}$  and  $-\frac{1}{2}N < [i-j]_N \leq \frac{1}{2}N$ . Note that the relation (2.1) holds only asymptotically as  $W \rightarrow \infty$  but this can be remedied by an irrelevant rescaling. If the bandwidth is comparable with  $N$ , then we also have to assume that  $f(x)$  is supported in  $|x| \leq N/(2W)$ .

The quantity  $M$  defined in (2.3) satisfies  $M \leq W/\|f\|_\infty$ . In Appendix A we will show that (2.5) is satisfied for the choice of (2.8) if  $W$  is large enough.

The Stieltjes transform of the empirical eigenvalue distribution of  $H$  is given by

$$m(z) \equiv m_N(z) = \frac{1}{N} \text{Tr} \frac{1}{H - z}, \quad z = E + i\eta. \quad (2.9)$$

We define the density of the semicircle law

$$\varrho_{sc}(x) := \frac{1}{2\pi} \sqrt{[4 - x^2]_+}, \quad (2.10)$$

and, for  $\text{Im } z > 0$ , its Stieltjes transform

$$m_{sc}(z) := \int_{\mathbb{R}} \frac{\varrho_{sc}(x)}{x - z} dx. \quad (2.11)$$

The Stieltjes transform  $m_{sc}(z) \equiv m_{sc}$  may also be characterized as the unique solution of

$$m_{sc} + \frac{1}{z + m_{sc}} = 0 \quad (2.12)$$

satisfying  $\text{Im } m_{sc}(z) > 0$  for  $\text{Im } z > 0$ , i.e.,

$$m_{sc}(z) = \frac{-z + \sqrt{z^2 - 4}}{2}. \quad (2.13)$$

Here the square root function is chosen with a branch cut along the positive real axis. This guarantees that the imaginary part of  $m_{sc}$  is non-negative. The Wigner semicircle law states that  $m_N(z) \rightarrow m_{sc}(z)$  for any fixed  $z$  provided that  $\eta = \text{Im } z > 0$  is independent of  $N$ . The local version of this result for universal Wigner matrices is the content of the following Theorem.

**Theorem 2.1 (Local semicircle law)** *Let  $H = (h_{ij})$  be a Hermitian  $N \times N$  random matrix where the matrix elements  $h_{ij} = \bar{h}_{ji}$ ,  $i \leq j$ , are independent random variables with  $\mathbb{E} h_{ij} = 0$ ,  $1 \leq i, j \leq N$ , and assume that the variances  $\sigma_{ij}^2 = \mathbb{E}|h_{ij}|^2$  satisfy (2.1), (2.4) and (2.5). Suppose that the distributions of the matrix elements have a uniformly subexponential decay in the sense that there exist constants  $\alpha, \beta > 0$ , independent of  $N$ , such that for any  $x > 0$  we have*

$$\mathbb{P}(|h_{ij}| \geq x^\alpha |\sigma_{ij}|) \leq \beta e^{-x}. \quad (2.14)$$

*Then there exist constants  $C_1, C_2, C$  and  $c > 0$ , depending only on  $\alpha, \beta$  and  $\delta_-$  in (2.5), such that for any  $z = E + i\eta$  with  $\eta = \text{Im } z > 0$ ,  $|z| \leq 10$  and*

$$\frac{1}{\sqrt{M\eta}} \leq \frac{\kappa^2}{(\log N)^{C_1}}, \quad (2.15)$$

*where  $\kappa := ||E| - 2|$ , the Stieltjes transform of the empirical eigenvalue distribution of  $H$  satisfies*

$$\mathbb{P} \left( |m_N(z) - m_{sc}(z)| \geq (\log N)^{C_2} \frac{1}{\sqrt{M\eta}\kappa} \right) \leq CN^{-c(\log \log N)} \quad (2.16)$$

for sufficiently large  $N$ . In fact, the same result holds for the individual matrix elements of the Green's function  $G_{ii}(z) = (H - z)^{-1}(i, i)$ :

$$\mathbb{P} \left( \max_i |G_{ii}(z) - m_{sc}(z)| \geq (\log N)^{C_2} \frac{1}{\sqrt{M\eta\kappa}} \right) \leq CN^{-c(\log \log N)}. \quad (2.17)$$

We remark that once a local semicircle law is obtained on a scale essentially  $M^{-1}$ , it is straightforward to show that eigenvectors are delocalized on a scale at least of order  $M$ . The precise statement will be formulated in Corollary 3.2. We will prove Theorem 2.1 in Sections 3–5 by extending the approach of [12, 13, 14]. The main ingredients of this approach consist of i) a derivation of a self-consistent equation for the Green's function and ii) an induction on the scale of the imaginary part of the energy. The key novelty in this paper is that the self-consistent equation is formulated for the array of the diagonal elements of the Green's function  $(G_{11}, G_{22}, \dots, G_{NN})$  instead of the Stieltjes transform  $m = \frac{1}{N} \text{Tr} G = \frac{1}{N} \sum_i G_{ii}$  itself as in [14]. This yields for the first time a strong pointwise control on the diagonal elements  $G_{ii}$ , see (2.17).

The subexponential decay condition (2.14) can be weakened if we are not aiming at error estimates faster than any power law of  $N$ . This can be easily carried out and we will not pursue it in this paper.

Denote the eigenvalues of  $H$  by  $\lambda_1, \dots, \lambda_N$  and let  $p_N(x_1, \dots, x_N)$  be their (symmetric) probability density. For any  $k = 1, 2, \dots, N$ , the  $k$ -point correlation function of the eigenvalues is defined by

$$p_N^{(k)}(x_1, x_2, \dots, x_k) := \int_{\mathbb{R}^{N-k}} p_N(x_1, x_2, \dots, x_N) dx_{k+1} \dots dx_N. \quad (2.18)$$

We now state our main result concerning these correlation functions.

**Theorem 2.2 (Universality for generalized Wigner matrices)** *We consider a generalized hermitian Wigner matrix such that (2.6) holds. Assume that the distributions  $\nu_{ij}$  of the  $(i, j)$  matrix elements have a uniformly subexponential decay in the sense of (2.14). Suppose that the real and imaginary parts of  $h_{ij}$  are i.i.d., distributed according to  $\omega_{ij}$ , i.e.,  $\nu_{ij}(dh) = \omega_{ij}(d\text{Im} h)\omega_{ij}(d\text{Re} h)$ . Let  $m_k(i, j) = \int x^k d\omega_{ij}(x)$ ,  $1 \leq k \leq 4$ , denote the  $k$ -th moment of  $\omega_{ij}$  ( $m_1 = 0$ ). Suppose that*

$$\inf_N \min_{1 \leq i, j \leq N} \left\{ \frac{m_4(i, j)}{(m_2(i, j))^2} - \frac{(m_3(i, j))^2}{(m_2(i, j))^3} \right\} > 1, \quad (2.19)$$

then, for any  $k \geq 1$  and for any compactly supported continuous test function  $O : \mathbb{R}^k \rightarrow \mathbb{R}$ , we have

$$\begin{aligned} & \lim_{b \rightarrow 0} \lim_{N \rightarrow \infty} \frac{1}{2b} \int_{E-b}^{E+b} dE' \int_{\mathbb{R}^k} d\alpha_1 \dots d\alpha_k O(\alpha_1, \dots, \alpha_k) \\ & \times \frac{1}{\varrho_{sc}(E)^k} \left( p_N^{(k)} - p_{GUE, N}^{(k)} \right) \left( E' + \frac{\alpha_1}{N\varrho_{sc}(E)}, \dots, E' + \frac{\alpha_k}{N\varrho_{sc}(E)} \right) = 0, \end{aligned} \quad (2.20)$$

where  $p_{GUE, N}^{(k)}$  is the  $k$ -point correlation function of the GUE ensemble. The same statement holds for generalized symmetric Wigner matrices, with GOE replacing the GUE ensemble.

The limiting correlation functions of the GUE ensemble are given by the sine kernel

$$\frac{1}{\varrho_{sc}(E)^k} p_{GUE, N}^{(k)} \left( E + \frac{\alpha_1}{N\varrho_{sc}(E)}, \dots, E + \frac{\alpha_k}{N\varrho_{sc}(E)} \right) \rightarrow \det \{ K(\alpha_i - \alpha_j) \}_{i, j=1}^k, \quad K(x) = \frac{\sin \pi x}{\pi x},$$

and similar universal formula is available for the limiting gap distribution.

*Remark:* The quantity in the bracket in (2.19) is always greater or equal to 1 for any real distribution with mean zero, which can be obtained by

$$m_3^2 = \left[ \int x^3 d\omega \right]^2 = \left[ \int x(x^2 - m_2) d\omega \right]^2 \leq \left[ \int x^2 d\omega \right] \left[ \int (x^2 - m_2)^2 d\omega \right] = m_2(m_4 - m_2^2)$$

and it is exactly 1 if the distribution is supported on two points. For example, if  $\omega_{ij}$  is a rescaling of a fixed distribution  $\tilde{\omega}$  with variance  $\frac{1}{2}$ , i.e.  $\omega_{ij}(x)dx = \sigma_{ij}^{-1}\tilde{\omega}(x/\sigma_{ij})dx$ , then condition (2.19) is satisfied under (2.6), as long as the support of  $\tilde{\omega}$  consists of at least three points. The case of a Bernoulli-type distribution supported on two points require a separate argument and it will be treated in the forthcoming paper [20].

We now state our main comparison theorem for matrix elements of Green's functions of two Wigner ensembles. As in the paper [31], we assume conditions on four moments. It will lead quickly to Theorem 6.4 stating that the correlation functions of eigenvalues of two matrix ensembles are identical up to scale  $1/N$  provided that the first four moments of all matrix elements of these two ensembles are almost identical. Here we do not assume that the real and imaginary parts are i.i.d., hence the  $k$ -th moment of  $h_{ij}$  is understood as the collection of numbers  $\int \tilde{h}^s h^{k-s} \nu_{ij}(dh)$ ,  $s = 0, 1, 2, \dots, k$ . The main result in [31] compares the joint distribution of individual eigenvalues — which is not covered by our Theorem 2.3 — but it does not address directly the matrix elements of Green's functions. The key input for both theorems is the local semicircle law on the almost optimal scale  $N^{-1+\varepsilon}$ . The eigenvalue perturbation used in [31] requires certain estimates on the eigenvalue level repulsion; the proof of Theorem 2.3 is a straightforward resolvent perturbation theory.

**Theorem 2.3 (Green's function comparison)** *Suppose that we have two generalized  $N \times N$  Wigner matrices,  $H^{(v)}$  and  $H^{(w)}$ , with matrix elements  $h_{ij}$  given by the random variables  $N^{-1/2}v_{ij}$  and  $N^{-1/2}w_{ij}$ , respectively, with  $v_{ij}$  and  $w_{ij}$  satisfying the uniform subexponential decay condition*

$$\mathbb{P}(|v_{ij}| \geq x^\alpha) \leq \beta e^{-x}, \quad \mathbb{P}(|w_{ij}| \geq x^\alpha) \leq \beta e^{-x},$$

with some  $\alpha, \beta > 0$ . Fix a bijective ordering map on the index set of the independent matrix elements,

$$\phi : \{(i, j) : 1 \leq i \leq j \leq N\} \rightarrow \{1, \dots, \gamma(N)\}, \quad \gamma(N) := \frac{N(N+1)}{2},$$

and denote by  $H_\gamma$  the generalized Wigner matrix whose matrix elements  $h_{ij}$  follow the  $v$ -distribution if  $\phi(i, j) \leq \gamma$  and they follow the  $w$ -distribution otherwise; in particular  $H^{(v)} = H_0$  and  $H^{(w)} = H_{\gamma(N)}$ . Let  $\kappa > 0$  be arbitrary and suppose that, for any small parameter  $\tau > 0$  and for any  $y \geq N^{-1+\tau}$ , we have the following estimate on the diagonal elements of the resolvent

$$\mathbb{P} \left( \max_{0 \leq \gamma \leq \gamma(N)} \max_{1 \leq k \leq N} \max_{|E| \leq 2^{-\kappa}} \left| \left( \frac{1}{H_\gamma - E - iy} \right)_{kk} \right| \leq N^{2\tau} \right) \geq 1 - CN^{-c \log \log N} \quad (2.21)$$

with some constants  $C, c$  depending only on  $\tau, \kappa$ . Moreover, we assume that the first three moments of  $v_{ij}$  and  $w_{ij}$  are the same, i.e.

$$\mathbb{E} \bar{v}_{ij}^s v_{ij}^u = \mathbb{E} \bar{w}_{ij}^s w_{ij}^u, \quad 0 \leq s + u \leq 3,$$

and the difference between the fourth moments of  $v_{ij}$  and  $w_{ij}$  is much less than 1, say

$$|\mathbb{E} \bar{v}_{ij}^s v_{ij}^{4-s} - \mathbb{E} \bar{w}_{ij}^s w_{ij}^{4-s}| \leq N^{-\delta}, \quad s = 0, 1, 2, 3, 4, \quad (2.22)$$



for some given  $\delta > 0$ . Let  $\varepsilon > 0$  be arbitrary and choose an  $\eta$  with  $N^{-1-\varepsilon} \leq \eta \leq N^{-1}$ . For any sequence of positive integers  $k_1, \dots, k_n$ , set complex parameters  $z_j^m = E_j^m \pm i\eta$ ,  $j = 1, \dots, k_m$ ,  $m = 1, \dots, n$ , with  $|E_j^m| \leq 2 - 2\kappa$  and with an arbitrary choice of the  $\pm$  signs. Let  $G^{(v)}(z) = (H^{(v)} - z)^{-1}$  denote the resolvent and let  $F(x_1, \dots, x_n)$  be a function such that for any multi-index  $\alpha = (\alpha_1, \dots, \alpha_n)$  with  $1 \leq |\alpha| \leq 5$  and for any  $\varepsilon' > 0$  sufficiently small, we have

$$\max \left\{ |\partial^\alpha F(x_1, \dots, x_n)| : \max_j |x_j| \leq N^{\varepsilon'} \right\} \leq N^{C_0 \varepsilon'} \quad (2.23)$$

and

$$\max \left\{ |\partial^\alpha F(x_1, \dots, x_n)| : \max_j |x_j| \leq N^2 \right\} \leq N^{C_0} \quad (2.24)$$

for some constant  $C_0$ .

Then, there is a constant  $C_1$ , depending on  $\alpha, \beta, \sum_m k_m$  and  $C_0$  such that for any  $\eta$  with  $N^{-1-\varepsilon} \leq \eta \leq N^{-1}$  and for any choices of the signs in the imaginary part of  $z_j^m$ , we have

$$\left| \mathbb{E} F \left( \frac{1}{N^{k_1}} \text{Tr} \left[ \prod_{j=1}^{k_1} G^{(v)}(z_j^1) \right], \dots, \frac{1}{N^{k_n}} \text{Tr} \left[ \prod_{j=1}^{k_n} G^{(v)}(z_j^n) \right] \right) - \mathbb{E} F \left( G^{(v)} \rightarrow G^{(w)} \right) \right| \leq C_1 N^{-1/2+C_1\varepsilon} + C_1 N^{-\delta+C_1\varepsilon}, \quad (2.25)$$

where the arguments of  $F$  in the second term are changed from the Green's functions of  $H^{(v)}$  to  $H^{(w)}$  and all other parameters remain unchanged.

*Remark 1:* We formulated Theorem 2.3 for functions of traces of monomials of the Green's function because this is the form we need in the application. However, the result (and the proof we are going to present) holds directly for matrix elements of monomials of Green's functions as well, namely, for any choice of  $\ell_1, \dots, \ell_{2n}$ , we have

$$\left| \mathbb{E} F \left( \frac{1}{N^{k_1-1}} \left[ \prod_{j=1}^{k_1} G^{(v)}(z_j^1) \right]_{\ell_1, \ell_2}, \dots, \frac{1}{N^{k_n-1}} \left[ \prod_{j=1}^{k_n} G^{(v)}(z_j^n) \right]_{\ell_{2n-1}, \ell_{2n}} \right) - \mathbb{E} F \left( G^{(v)} \rightarrow G^{(w)} \right) \right| \leq C_1 N^{-1/2+C_1\varepsilon} + C_1 N^{-\delta+C_1\varepsilon}. \quad (2.26)$$

We also remark that Theorem 2.3 holds for generalized Wigner matrices since  $C_{sup} < \infty$  in (2.2). The positive lower bound on the variances,  $C_{inf} > 0$ , is not necessary for this theorem.

*Remark 2:* Although we state Theorem 2.3 for Hermitian and symmetric ensembles, similar results hold for real and complex sample covariance ensembles; the modification of the proof, to be given in Section 8, is obvious and we omit the details.

To summarize, our approach to prove the universality is based on the following three steps; a detailed outline will be given in Section 6. *Step 1.* Local semicircle law, i.e., Theorem 2.1. This will be proved in Sections 3–5. *Step 2.* Universality for ensembles with smooth distributions satisfying the logarithmic Sobolev inequality (LSI), Theorem 6.3. The key input is the general theorem, Theorem 6.2, concerning the universality for the local relaxation flow. In Section 7, by using the local semicircle law and the LSI, we

verify the assumptions for this theorem. *Step 3.* Green's function comparison theorem, Theorem 2.3. This removes the restriction on the smoothness and the LSI, and it will be proved in Section 8.

*Convention.* We will frequently use the notation  $C, c$  for generic positive constants whose exact values are irrelevant and may change from line to line. For two positive quantities  $A, B$  we also introduce the notation  $A \asymp B$  to indicate that there exists a universal constant  $C$  such that  $C^{-1} \leq A/B \leq C$ .

### 3 Proof of local semicircle law

*Proof of Theorem 2.1* Recall that  $G_{ij} = G_{ij}(z)$  denotes the matrix element

$$G_{ij} = \left( \frac{1}{H - z} \right)_{ij} \quad (3.1)$$

and

$$m(z) = m_N(z) = \sum_{i=1}^N N^{-1} G_{ii}(z).$$

We will prove the following more detailed stronger result.

**Theorem 3.1** *Assume the  $N \times N$  random matrix  $H$  satisfies (2.1), (2.4), (2.5) and (2.14),  $\mathbb{E} h_{ij} = 0$ , for any  $1 \leq i, j \leq N$ . Let  $z = E + i\eta$  ( $\eta > 0$ ) and let  $g(z)$  be the real valued function defined by*

$$g(z) \equiv \min \left\{ \sqrt{\kappa + \eta}, \max \{ \delta_+, |\operatorname{Re} [m_{sc}(z)^2] - 1| \} \right\}, \quad (3.2)$$

where  $\kappa \equiv ||E| - 2|$  and  $\delta_+$  is given in (2.4). Then for all  $z = E + i\eta$  and

$$\frac{(\kappa + \eta)^{1/4}}{\sqrt{M\eta} g^2(z)} \leq (\log N)^{-13-6\alpha}, \quad |z| \leq 10, \quad (3.3)$$

we have

$$\mathbb{P} \left\{ \max_i |G_{ii}(z) - m_{sc}(z)| \geq (\log N)^{11+6\alpha} \frac{(\kappa + \eta)^{1/4}}{\sqrt{M\eta} g(z)} \right\} \leq CN^{-c(\log \log N)} \quad (3.4)$$

for sufficiently large  $N$ , with positive  $c$  and  $C > 0$  depending only  $\alpha$  and  $\beta$  in (2.14) and  $\delta_-$  in (2.4) and (2.5).

*Remark:* The condition (3.3) is effectively a lower bound on  $\eta$ . The control function  $g(z)$  can be estimated by

$$g(z) \asymp \begin{cases} \min \left\{ \sqrt{\kappa + \eta}, \max \{ \delta_+, \eta/\sqrt{\kappa}, \kappa \} \right\}, & |E| \leq 2 \text{ and } \kappa \geq \eta, \\ \sqrt{\kappa + \eta}, & \text{otherwise,} \end{cases} \quad (3.5)$$

up to some factor of order one. Note that the precise formula (3.2) for  $g(z)$  is not important, only its asymptotic behaviour for small  $\kappa, \eta$  and  $\delta_+$  is relevant. The theorem remains valid if  $g(z)$  is replaced by  $\tilde{g}(z)$  with  $\tilde{g}(z) \leq Cg(z)$ . In particular,  $g(z)$  can be chosen to be order one when  $E$  is not near the edges of the spectrum. If we are only concerned with the case of generalized Wigner matrices, (2.6), we can choose

$g(z) = O(\sqrt{\kappa + \eta})$  for any  $z = E + i\eta$  ( $\eta > 0$ ). Note that Theorem 2.1 was obtained by replacing  $g(z)$  with the lower bound  $\kappa \leq g(z)$  in Theorem 3.1.

Once the local semicircle law is established on scale  $\eta \asymp 1/M$  (modulo logarithmic factors), we obtain the following supremum bound on the eigenvectors that can be interpreted as a lower bound of order  $1/M$  on the localization length. The proof of this result now is simpler than in [12, 13], since we have a pointwise control on the diagonal elements of the Green's function. Let  $\mathbf{u}_\alpha$  denote the normalized eigenvector of  $H$  belonging to the eigenvalue  $\lambda_\alpha$ ,  $\alpha = 1, 2, \dots, N$ , i.e.,  $H\mathbf{u}_\alpha = \lambda_\alpha\mathbf{u}_\alpha$  and  $\|\mathbf{u}_\alpha\| = 1$ .

**Corollary 3.2** *Let  $H$  be as in Theorem 3.1, for any fixed  $\kappa > 0$ , there exists  $C_\kappa$  that*

$$\mathbb{P} \left\{ \exists \lambda_\alpha \in [-2 + \kappa, 2 - \kappa], H\mathbf{u}_\alpha = \lambda_\alpha\mathbf{u}_\alpha, \|\mathbf{u}_\alpha\| = 1, \|\mathbf{u}_\alpha\|_\infty \geq C_\kappa \frac{(\log N)^{13+6\alpha}}{M^{1/2}} \right\} \leq CN^{-c \log \log N}. \quad (3.6)$$

For the case of generalized Wigner matrices, (2.6), we have the following more precise bound

$$\mathbb{P} \left\{ \exists \lambda_\alpha, H\mathbf{u}_\alpha = \lambda_\alpha\mathbf{u}_\alpha, \|\mathbf{u}_\alpha\| = 1, \|\mathbf{u}_\alpha\|_\infty \geq \frac{C(\log N)^{13+6\alpha}}{N^{1/2} [|\lambda_\alpha| - 2| + N^{-1}]^{1/2}} \right\} \leq CN^{-c \log \log N}. \quad (3.7)$$

*Proof of Corollary 3.2.* Let  $\eta = C_\kappa M^{-1}(\log N)^{26+12\alpha}$ ;  $C_\kappa$  can be chosen large enough so that (3.3) is satisfied for all  $|\kappa'| \leq \kappa$ , making use of (3.5). Choose  $\{E_m\}$  as a grid of points in  $[-2 + \kappa, 2 - \kappa]$  such that the distance between any two neighbors is of order  $\eta$ . Then with (3.4), we have

$$\mathbb{P} \left\{ \max_j \max_m \operatorname{Im} |G_{jj}(E_m + i\eta)| \geq \operatorname{Im} m_{sc}(E_m + i\eta) + 1 \right\} \leq CN^{-c(\log \log N)}, \quad (3.8)$$

where we used  $g(z) \leq \sqrt{\kappa + \eta} \leq C$  from (3.5). Then, with  $|m_{sc}(z)| \leq C$  (see (2.13)) and

$$\operatorname{Im} G_{jj}(E_m + i\eta) = \sum_\alpha \frac{\eta |u_\alpha(j)|^2}{|E_m - \lambda_\alpha|^2 + \eta^2}, \quad (3.9)$$

where  $\mathbf{u}_\alpha = (u_\alpha(1), u_\alpha(2), \dots, u_\alpha(N))$ , we have

$$\mathbb{P} \left( \max_j \max_m \sum_\alpha \frac{\eta |u_\alpha(j)|^2}{|E_m - \lambda_\alpha|^2 + \eta^2} \geq C \right) \leq CN^{-c(\log \log N)}. \quad (3.10)$$

By the definition of  $E_m$ , for any  $\lambda_\alpha \in [-2 + \kappa, 2 - \kappa]$ , there exists  $m'$  such that  $|E_{m'} - \lambda_\alpha|$  is of the order of  $\eta$ . Together with (3.10), we obtain (3.6).

In case of the generalized Wigner matrix (2.6), we have  $g(z) = \sqrt{\kappa + \eta}$  and  $M \asymp N$ . Let  $\eta$  be the solution to  $\eta = N^{-1}(\log N)^{26+12\alpha}(\kappa + \eta)^{-3/2}$ , then  $N^{-1} \leq \eta \leq CN^{-1}(\kappa + N^{-1})^{-3/2}(\log N)^{26+12\alpha}$ . With this choice of  $\eta$ , (3.3) is satisfied, and  $\max_i |G_{ii} - m_{sc}| \leq C(\log N)^{-2}(\kappa + \eta)^{1/2}$  holds with an overwhelming probability by (3.4). Since  $|\operatorname{Im} m_{sc}(z)| \leq C\sqrt{\kappa + \eta}$ , so  $\max_i \operatorname{Im} G_{ii} \leq C(\kappa + \eta)^{1/2}$ . By the argument above, we obtain that  $\|\mathbf{u}_\alpha\|_\infty^2 \leq C\eta(\kappa + \eta)^{1/2}$  on this event. This proves (3.7).  $\square$

To prove that  $G_{ii}(z)$  is very close to  $m_{sc}(z)$  in the sense of (3.4), we will also need to control the off-diagonal elements. In fact we will show that all  $G_{ij}$  ( $i \neq j$ ) are bounded by  $O((M\eta)^{-1})$  up to some factor  $(\log N)^C$ . To state the result precisely, we first define some events in the probability space.

Recall that  $\lambda_\alpha$ ,  $\alpha = 1, 2, \dots, N$ , denote the eigenvalues of  $H = (h_{ij})$ . Denote by  $\Omega^0$  the subset of the probability space such that

$$\max_\alpha |\lambda_\alpha| \leq 3. \quad (3.11)$$

Let  $\widehat{\Omega}_z^d$  (here the superscript  $d$  means diagonal) be the subset of  $\Omega^0$  where the following inequality on the diagonal terms hold for any  $1 \leq i \leq N$

$$|G_{ii}(z) - m_{sc}(z)| \leq (\log N)^{11+6\alpha} \frac{(\kappa + \eta)^{1/4}}{\sqrt{M\eta} g(z)} \quad (3.12)$$

(recall that  $m_{sc}(z)$  was defined in (2.13)). Similarly, let  $\widehat{\Omega}_z^o$  (here the superscript  $o$  means off-diagonal) be the subset of  $\Omega^0$  where the following inequality on the off-diagonal terms hold for any  $1 \leq i \neq j \leq N$

$$|G_{ij}(z)| \leq (\log N)^{5+4\alpha} \frac{(\kappa + \eta)^{1/4}}{\sqrt{M\eta}}. \quad (3.13)$$

Finally, denote by  $\Omega_z^d$  the set

$$\Omega_z^d = \bigcap_{k=0}^{10N^5} \widehat{\Omega}_{z+ik/N^5}^d, \quad (3.14)$$

and similarly define  $\Omega_z^o$ . These sets depend on  $N$  but we suppress this from the notations.

*Proof of Theorem 3.1.* The following proposition immediately implies Theorem 3.1. □

**Proposition 3.3** *Suppose that the assumptions of Theorem 3.1 hold. Then, for sufficiently large  $N$ , we have*

$$\mathbb{P}(\Omega_z^d \cap \Omega_z^o) \geq 1 - CN^{-c \log \log N} \quad (3.15)$$

for some positive constants  $c$  and  $C$ .

Following the work of [12], we will use a continuity argument. In Section 4 we will derive a self-consistent equation of the form

$$G_{ii} + \frac{1}{z + \sum_j \sigma_{ij}^2 G_{jj} + \Upsilon_i(z)} = 0, \quad i = 1, 2, \dots, N. \quad (3.16)$$

Later we will give an explicit formula for  $\Upsilon_i(z)$ , but for now we take (3.16) as the definition of  $\Upsilon_i$ . Let  $\widehat{\Omega}_z^\Upsilon(N) = \widehat{\Omega}_z^\Upsilon$  be the subset of  $\Omega^0$  where the following inequality holds

$$\Upsilon = \Upsilon(z) := \max_i |\Upsilon_i(z)| \leq (\log N)^{9+6\alpha} \frac{(\kappa + \eta)^{1/4}}{\sqrt{M\eta}}. \quad (3.17)$$

We will use the following Lemmas that will be proved later in Section 4 and 5.

**Lemma 3.4** *Let  $z = E + i\eta$  be a fixed complex number satisfying (3.3). Then there are constants  $C$  and  $c$  such that for  $N \geq N_0$ , with  $N_0$  sufficiently large independent of  $E$  and  $\eta$ , the following estimates hold.*

(1) *Suppose  $3 \leq \eta \leq 10$ . Then*

$$\mathbb{P}(\widehat{\Omega}_z^o) \geq 1 - CN^{-c \log \log N}, \quad (3.18)$$

$$\mathbb{P}(\widehat{\Omega}_z^o \cap \widehat{\Omega}_z^\Upsilon) \geq 1 - CN^{-c \log \log N}, \quad (3.19)$$

and for any  $1 \leq i \leq N$ ,

$$|\mathbb{E}\mathbf{1}(\widehat{\Omega}_z^o)\Upsilon_i(z)| \leq (\log N)^{10+8\alpha} \frac{(\kappa + \eta)^{1/2}}{M\eta} + CN^{-c(\log \log N)}. \quad (3.20)$$

(2) *Suppose that  $\eta \leq 3$ . Setting  $z' = z + iN^{-5}$ , we have*

$$\mathbb{P}(\widehat{\Omega}_z^o \cap \Omega_{z'}^d \cap \Omega_{z'}^o) \geq \mathbb{P}(\Omega_{z'}^d \cap \Omega_{z'}^o) - CN^{-c \log \log N}, \quad (3.21)$$

$$\mathbb{P}(\widehat{\Omega}_z^o \cap \Omega_{z'}^d \cap \Omega_{z'}^o \cap \widehat{\Omega}_z^\Upsilon) \geq \mathbb{P}(\widehat{\Omega}_z^o \cap \Omega_{z'}^d \cap \Omega_{z'}^o) - CN^{-c \log \log N}, \quad (3.22)$$

and for any  $1 \leq i \leq N$ ,

$$|\mathbb{E}\mathbf{1}(\widehat{\Omega}_z^o \cap \Omega_{z'}^d \cap \Omega_{z'}^o)\Upsilon_i(z)| \leq (\log N)^{10+8\alpha} \frac{(\kappa + \eta)^{1/2}}{M\eta} + CN^{-c(\log \log N)}. \quad (3.23)$$

**Lemma 3.5** *Suppose we are on the event  $\widehat{\Omega}_z^\Upsilon$  for some fixed  $z = E + i\eta$  satisfying (3.3). Suppose either  $3 \leq \eta \leq 10$  or the following inequality hold:*

$$\max_i |G_{ii} - m_{sc}(z)| \leq 2(\log N)^{-2}g(z). \quad (3.24)$$

Then, for sufficiently large  $N$ , we have

$$\max_i |G_{ii} - m_{sc}(z)| \leq \frac{(\log N)^2}{g(z)} \Upsilon(z). \quad (3.25)$$

*Proof of Proposition 3.3.* Recall that  $\widehat{\Omega}_z^d$  is the subset of  $\Omega^0$  where (3.12) holds. Since on  $\widehat{\Omega}_z^\Upsilon$  (3.12) follows from (3.17), the case  $3 \leq \eta = \text{Im } z \leq 10$  follows from Lemma 3.4 and Lemma 3.5 by taking a union bound for  $0 \leq k \leq 10N^5$ .

Now we prove (3.15) for the case  $\eta \leq 3$  assuming that  $z = E + i\eta$  satisfies (3.3). We have shown that (3.15) holds for  $\eta = 3$ , now we will successively decrease  $\eta$  by  $N^{-5}$  in each step, and we continue this inductive procedure as long as (3.3) is still satisfied for the reduced  $\eta$ . More precisely, let  $z' = z + iN^{-5}$  and assume that (3.15) holds for  $z'$ . Our goal is to prove that

$$\mathbb{P}(\Omega_z^d \cap \Omega_z^o) \geq \mathbb{P}(\Omega_{z'}^d \cap \Omega_{z'}^o) - CN^{-c \log \log N}. \quad (3.26)$$

The number of steps we will be taking is of order  $N^5$ . Since  $N^{-c \log \log N} \ll N^{-5}$ , this proves (3.15) provided that we can establish (3.26).

From (3.21), the difference between the probabilities of the sets  $\widehat{\Omega}_z^o \cap \Omega_{z'}^d, \cap \Omega_z^o$ , and  $\Omega_{z'}^d \cap \Omega_z^o$ , is negligible. With the definition of  $\Omega_z^d$  and  $\Omega_z^o$  in (3.14), we have

$$\Omega_z^d \cap \Omega_z^o \supset \widehat{\Omega}_z^d \cap \widehat{\Omega}_z^o \cap \Omega_{z'}^d \cap \Omega_{z'}^o. \quad (3.27)$$

Then, to prove (3.26), it remains to prove

$$\mathbb{P}(\widehat{\Omega}_z^d \cap \widehat{\Omega}_z^o \cap \Omega_{z'}^d \cap \Omega_{z'}^o) \geq \mathbb{P}(\widehat{\Omega}_z^o \cap \Omega_{z'}^d \cap \Omega_{z'}^o) - CN^{-c \log \log N}, \quad (3.28)$$

i.e., we need to estimate the probability of the complement of  $\widehat{\Omega}_z^d$  on the set  $\widehat{\Omega}_z^o \cap \Omega_{z'}^d \cap \Omega_{z'}^o$ . On this set, using (3.22), we can assume that the estimate (3.17) holds with a very high probability. We will show below that (3.24) holds on  $\Omega_{z'}^d$ . Then (3.25) together with (3.17) imply (3.12), the defining relation of  $\widehat{\Omega}_z^d$ . This will conclude (3.28) and complete the proof of Proposition 3.3. Therefore, we only have to verify (3.24).

Now we show that (3.24) holds on  $\Omega_{z'}^d$ . Recall  $z' = z + iN^{-5}$  and we have the trivial estimate

$$|G_{ii}(z) - m_{sc}(z)| \leq |G_{ii}(z) - G_{ii}(z')| + |m_{sc}(z) - m_{sc}(z')| + |G_{ii}(z') - m_{sc}(z')|. \quad (3.29)$$

In the set  $\Omega_{z'}^d$ , we have

$$|G_{ii}(z') - m_{sc}(z')| \leq (\log N)^{11+6\alpha} \frac{(\kappa + \eta)^{1/4}}{\sqrt{M\eta} g(z')} \leq (\log N)^{-2} g(z'), \quad (3.30)$$

where in the second inequality we used (3.3). By the definition of  $g(z)$  from (3.2), we have  $g(z) \leq \sqrt{\kappa + \eta}$ . Thus, if (3.3) holds, then, in particular,

$$\eta \geq CM^{-1}(\log N)^{26+12\alpha}. \quad (3.31)$$

This sets a lower bound on  $\eta$ . Together with  $|z - z'| = 1/N^5$ , we have the trivial continuity bound

$$|G_{ii}(z) - G_{ii}(z')| + |m_{sc}(z) - m_{sc}(z')| \leq N^{-2},$$

using  $|\partial_z m_{sc}(z)| \leq |\text{Im } z|^{-2}$ ,  $|\partial_z G_{ii}(z)| = |[(H-z)^{-2}]_{ii}| \leq \|(H-z)^{-2}\| \leq |\text{Im } z|^{-2}$  and  $\eta > N^{-1}$  from (3.31). Thus

$$|G_{ii}(z) - m_{sc}(z)| \leq N^{-2} + (\log N)^{-2} g(z'). \quad (3.32)$$

Using  $|g(z)| \geq C\eta \geq CN^{-1}$  and  $|g'(z)| \leq C\eta^{-1} \leq CN$  for  $\eta \leq 3$ , we have the following estimate

$$|G_{ii}(z) - m_{sc}(z)| \leq 2(\log N)^{-2} g(z) \quad (3.33)$$

in the set  $\Omega_{z'}^d$ . Thus the assumption (3.24) holds in the set  $\Omega_{z'}^d$ .  $\square$

Under the assumptions of Theorem 3.1, with (3.15), (3.20), (3.23) and the definitions in (3.14), all these  $\Omega$ 's are sets of almost full probability, i.e.,

$$\mathbb{P}(\widehat{\Omega}_z^o), \mathbb{P}(\widehat{\Omega}_z^d), \mathbb{P}(\Omega_z^o), \mathbb{P}(\Omega_z^d), \mathbb{P}(\widehat{\Omega}_z^Y) \geq 1 - CN^{c \log \log N} \quad (3.34)$$

for some  $c, C > 0$ .

## 4 Self-consistent equation for Green's function

First, we introduce some notations.

**Definition 4.1** For any collection of  $s$  different numbers,  $k_1, k_2, \dots, k_s \in \{1, 2, \dots, N\}$ , let  $H^{(k_1, k_2, \dots, k_s)}$  denote the  $N-s$  by  $N-s$  submatrix of  $H$  after removing the  $k_i$ -th ( $1 \leq i \leq s$ ) rows and columns. Sometimes we use the notation  $H^{(\mathbf{T})}$  where  $\mathbf{T}$  denote the unordered set  $\{k_1, k_2, \dots, k_s\}$ . Similarly, we define  $\mathbf{a}^{(\ell; \mathbf{T})}$  to be the  $\ell$ -th column of  $H$  with  $k_i$ -th ( $1 \leq i \leq s$ ) elements removed. Sometimes, we just use the short notation  $\mathbf{a}^\ell = \mathbf{a}^{(\ell; \mathbf{T})}$ .

For  $\mathbf{T} = \{k_1, k_2, \dots, k_s\}$ , we define

$$\begin{aligned} G_{ij}^{(\mathbf{T})} &:= [H^{(\mathbf{T})} - z]^{-1}(i, j), \\ Z_{ij}^{(\mathbf{T})} &:= \mathbf{a}^i \cdot [H^{(\mathbf{T})} - z]^{-1} \mathbf{a}^j = \sum_{k, l \notin \mathbf{T}} \overline{\mathbf{a}}_k^i G_{k, l}^{(\mathbf{T})} \mathbf{a}_l^j, \\ K_{ij}^{(\mathbf{T})} &:= h_{ij} - z\delta_{ij} - Z_{ij}^{(\mathbf{T})}. \end{aligned}$$

These quantities depend on  $z$ , but we mostly neglect this dependence in the notation.

We start the proof with deriving some identities between the matrix elements of  $G = (H - z)^{-1}$  and  $G^{(k_1, k_2, \dots, k_s)}$  using the following well known result in linear algebra that we quote without proof.

**Lemma 4.1** Let  $A$ ,  $B$ ,  $C$  be  $n \times n$ ,  $m \times n$  and  $m \times m$  matrices. We define  $(m+n) \times (m+n)$  matrix  $D$  as

$$D = \begin{pmatrix} A & B^* \\ B & C \end{pmatrix} \quad (4.1)$$

and  $n \times n$  matrix  $\widehat{D}$  as

$$\widehat{D} = A - B^* C^{-1} B. \quad (4.2)$$

Then for any  $1 \leq i, j \leq n$ , we have

$$(D^{-1})_{ij} = (\widehat{D}^{-1})_{ij}$$

for the corresponding matrix elements.

Furthermore, let  $\mathbf{T}$  denote the unordered set  $\{k_1, k_2, \dots, k_s\}$  and  $1 \leq k_i \leq n$ ,  $1 \leq i \leq s$ . We define  $D^{(\mathbf{T})}$  to be the  $n+m-s$  by  $n+m-s$  submatrix of  $D$  after removing the  $k_i$ -th ( $1 \leq i \leq s$ ) rows and columns and define  $\widehat{D}^{(\mathbf{T})}$  to be the  $n-s$  by  $n-s$  submatrix of  $\widehat{D}$  after removing the  $k_i$ -th ( $1 \leq i \leq s$ ) rows and columns. Then for any  $1 \leq i, j \leq n$  and  $i, j \notin \mathbf{T}$ , we have

$$\left( (D^{(\mathbf{T})})^{-1} \right)_{ij} = \left( (\widehat{D}^{(\mathbf{T})})^{-1} \right)_{ij}$$

for the corresponding matrix elements.

Using Lemma 4.1 and Definition 4.1, for  $1 \leq i \neq j \leq N$ , we have

$$G_{ii} = (K_{ii}^{(i)})^{-1} = \frac{K_{jj}^{(ij)}}{K_{jj}^{(ij)} K_{ii}^{(ij)} - K_{ij}^{(ij)} K_{ji}^{(ij)}}. \quad (4.3)$$

For the off diagonal matrix elements  $G_{ij}$ , ( $i \neq j$ ), we have

$$G_{ij} = -\frac{K_{ij}^{(ij)}}{K_{jj}^{(ij)} K_{ii}^{(ij)} - K_{ij}^{(ij)} K_{ji}^{(ij)}} = -G_{ii} \frac{K_{ij}^{(ij)}}{K_{jj}^{(ij)}} = -G_{ii} G_{jj}^{(i)} K_{ij}^{(ij)}. \quad (4.4)$$

Similarly, we have the following result

**Lemma 4.2** *Let  $\mathbf{T}$  be an unordered set  $\{k_1, k_2, \dots, k_s\}$  with  $1 \leq k_t \leq N$  for  $(1 \leq t \leq s)$  or  $\mathbf{T} = \emptyset$ . For simplicity, we use the notation  $(i \mathbf{T})$  for  $\{i\} \cup \mathbf{T}$  and  $(ij \mathbf{T})$  for  $\{i, j\} \cup \mathbf{T}$ . Then we have the following identities:*

1. For any  $i \notin \mathbf{T}$

$$G_{ii}^{(\mathbf{T})} = (K_{ii}^{(i \mathbf{T})})^{-1}. \quad (4.5)$$

2. For  $i \neq j$  and  $i, j \notin \mathbf{T}$

$$G_{ij}^{(\mathbf{T})} = -G_{jj}^{(\mathbf{T})} G_{ii}^{(j \mathbf{T})} K_{ij}^{(ij \mathbf{T})} = -G_{ii}^{(\mathbf{T})} G_{jj}^{(i \mathbf{T})} K_{ij}^{(ij \mathbf{T})}. \quad (4.6)$$

3. For  $i \neq j$  and  $i, j \notin \mathbf{T}$

$$G_{ii}^{(\mathbf{T})} - G_{ii}^{(j \mathbf{T})} = G_{ij}^{(\mathbf{T})} G_{ji}^{(\mathbf{T})} (G_{jj}^{(\mathbf{T})})^{-1}. \quad (4.7)$$

4. For any indices  $i, j$  and  $k$  that are different and  $i, j, k \notin \mathbf{T}$

$$G_{ij}^{(\mathbf{T})} - G_{ij}^{(k \mathbf{T})} = G_{ik}^{(\mathbf{T})} G_{kj}^{(\mathbf{T})} (G_{kk}^{(\mathbf{T})})^{-1}. \quad (4.8)$$

*Proof of Lemma 4.2.* The first two identities (4.5) and (4.6) are obvious extensions of (4.3) and (4.4). To prove (4.7), without loss of generality, we may assume that  $i = 1, j = 2$  and  $\mathbf{T} = \emptyset$ . Let  $D = H - z$  and  $\widehat{D}$  defined as in (4.2) with  $n = 2$  and  $m = N - 2$ . With Lemma 4.1, we have that for  $i, j = 1$  or  $2$ .

$$G_{ij} = (\widehat{D}^{-1})_{ij} \quad (4.9)$$

and

$$G_{ii}^{(j)} = \left( (\widehat{D}^{(j)})^{-1} \right)_{ii}. \quad (4.10)$$

Since  $\widehat{D}$  is just a  $2 \times 2$  matrix, one can easily check that (4.7) holds. With the same method, one can obtain (4.8).  $\square$

**Lemma 4.3** *The diagonal matrix elements of the resolvent satisfy the following self-consistent equation.*

$$G_{ii} = \left( -z - \sum_j \sigma_{ij}^2 G_{jj} + \Upsilon_i \right)^{-1} \quad (4.11)$$

where  $\Upsilon_i(z)$  is given by

$$\Upsilon_i(z) := \sigma_{ii}^2 G_{ii} + \sum_{j \neq i} \sigma_{ij}^2 G_{ij} G_{ji} [G_{ii}]^{-1} + \left( K_{ii}^{(i)} - \mathbb{E}_{\mathbf{a}^i} K_{ii}^{(i)} \right), \quad (4.12)$$



and  $\mathbb{E}_{\mathbf{a}^i}$  is the expectation over  $\mathbf{a}^i$ . Let  $\mathbf{T}$  denote the set  $\{k_1, k_2, \dots, k_m\}$ , which also could be the empty set, then

$$|K_{ii}^{(i\mathbf{T})} - \mathbb{E}_{\mathbf{a}^i} K_{ii}^{(i\mathbf{T})}| \leq (\log N)^{3+2\alpha} \sqrt{M^{-1} + M^{-1} \max_k |G_{kk}^{(i\mathbf{T})}|^2 + \max_{k \neq l} |G_{kl}^{(i\mathbf{T})}|^2} \quad (4.13)$$

and for  $i \neq j$

$$|K_{ij}^{(ij\mathbf{T})}| \leq (\log N)^{4+4\alpha} \sqrt{M^{-1} + (M\eta)^{-1} \max_l \left\{ \text{Im} G_{ll}^{(ij\mathbf{T})} \right\}}, \quad i \neq j, \quad (4.14)$$

hold with a probability larger than  $1 - CN^{-c(\log \log N)}$  for sufficiently large  $N$ .

*Proof of Lemma 4.3.* We can write  $G_{11}$  as follows,

$$G_{11} = (K_{11}^{(1)})^{-1} = \frac{1}{\mathbb{E}_{\mathbf{a}^1} K_{11}^{(1)} + K_{11}^{(1)} - \mathbb{E}_{\mathbf{a}^1} K_{11}^{(1)}}. \quad (4.15)$$

Using the fact  $G^{(1)} = (H^{(1)} - z)^{-1}$  is independent of  $\mathbf{a}^1$  and  $\mathbb{E}_{\mathbf{a}^1} \overline{\mathbf{a}^1(i)} \mathbf{a}^1(j) = \delta_{ij} \sigma_{1j}^2$ , we obtain  $\mathbb{E}_{\mathbf{a}^1} K_{11}^{(1)} = -z - \sum_{j \neq 1} \sigma_{1j}^2 G_{jj}^{(1)}$ , and thus

$$G_{11} = \frac{1}{-z - \sum_{j \neq 1} \sigma_{1j}^2 G_{jj}^{(1)} + (K_{11}^{(1)} - \mathbb{E}_{\mathbf{a}^1} K_{11}^{(1)})}. \quad (4.16)$$

Combining this identity with (4.7), we have

$$G_{11} = \left( -z - \sum_{j \neq 1} \sigma_{1j}^2 (G_{jj} - G_{1j} G_{j1} G_{11}^{-1}) + (K_{11}^{(1)} - \mathbb{E}_{\mathbf{a}^1} K_{11}^{(1)}) \right)^{-1}. \quad (4.17)$$

Clearly  $G_{11}$  can be replaced with any  $G_{ii}$  and this proves (4.11) with the definition (4.12).

Now we prove (4.13) and (4.14). Define

$$v_{ij} \equiv h_{ij} / \sigma_{ij}, \quad (4.18)$$

hence  $\mathbb{E} v_{ij} = 0$  and  $\mathbb{E} |v_{ij}|^2 = 1$ . If  $\sigma_{ij} = 0$ , i.e.,  $h_{ij} = 0$  almost surely, then we set  $v_{ij} = 0$ . By the definition of  $K_{ii}^{(i\mathbf{T})}$ , we write

$$K_{ii}^{(i\mathbf{T})} = h_{ii} - z - \sum_{k, l \notin (i\mathbf{T})} \overline{\mathbf{a}_k^i} G_{kl}^{(i\mathbf{T})} \mathbf{a}_l^i = h_{ii} - z - \sum_{k, l \notin (i\mathbf{T})} \overline{v_{ik}} \sigma_{ik} G_{kl}^{(i\mathbf{T})} \sigma_{li} v_{li} \quad (4.19)$$

and

$$\mathbb{E}_{\mathbf{a}^i} K_{ii}^{(i\mathbf{T})} = -z - \sum_{k \notin (i\mathbf{T})} \sigma_{ik} G_{kk}^{(i\mathbf{T})} \sigma_{ki}. \quad (4.20)$$

We note  $h_{ii}$ ,  $v_{ij}$  and  $G_{kl}^{(i\mathbf{T})}$  are independent for  $k, l \notin (i\mathbf{T})$ . With the sub-exponential decay (2.14) and  $\sigma_{ij}^2 \leq 1/M$ , we have for any  $i, j$

$$\mathbb{P} \left\{ |h_{ij}| \leq (\log N)^{3+2\alpha} M^{-1/2} \right\} \geq 1 - CN^{-c(\log \log N)}. \quad (4.21)$$

In Corollary B.3 of Appendix B we will prove a general large deviations result. Applying (B.15) to the last term in (4.19), with the choice

$$B_{kl} = \sigma_{ik} G_{kl}^{(i\mathbf{T})} \sigma_{li} \quad (4.22)$$

and with  $\sum_j \sigma_{ij}^2 = 1$  and  $\sigma_{ii}^2 \leq 1/M$ , we obtain that

$$\left| \sum_{k,l \notin (i\mathbf{T})} \overline{v_{ik}} \sigma_{ik} G_{kl}^{(i\mathbf{T})} \sigma_{li} v_{li} - \sum_{k \notin (i\mathbf{T})} \sigma_{ik} G_{kk}^{(i\mathbf{T})} \sigma_{ki} \right| \leq (\log N)^{3+2\alpha} \sqrt{M^{-1} \max_k |G_{kk}^{(i\mathbf{T})}|^2 + \max_{k \neq l} |G_{kl}^{(i\mathbf{T})}|^2} \quad (4.23)$$

holds with a probability larger than  $1 - CN^{-c(\log \log N)}$ . Together with (4.21), we obtain that (4.13) holds with a probability larger than  $1 - CN^{-c(\log \log N)}$  for sufficiently large  $N$ .

Next we prove (4.14). By the definition of  $K_{ij}^{(ij\mathbf{T})}$ ,  $i \neq j$ , we can write

$$K_{ij}^{(ij\mathbf{T})} = h_{ij} - \sum_{k,l \notin (ij\mathbf{T})} \overline{v_{ik}} \sigma_{ik} G_{kl}^{(ij\mathbf{T})} \sigma_{lj} v_{lj}. \quad (4.24)$$

Applying (B.16), (4.21) and  $\sigma_{ij}^2 \leq 1/M$ , we obtain that

$$|K_{ij}^{(ij\mathbf{T})}| \leq (\log N)^{4+4\alpha} \sqrt{M^{-1} + \sum_{k,l \notin (ij\mathbf{T})} |\sigma_{ik} G_{kl}^{(ij\mathbf{T})} \sigma_{lj}|^2} \quad (4.25)$$

holds with a probability larger than  $1 - CN^{-c(\log \log N)}$  for sufficiently large  $N$ . With Schwarz's inequality, for any  $i, j$ ,

$$\sum_{kl} |\sigma_{ik} G_{kl}^{(ij\mathbf{T})} \sigma_{lj}|^2 \leq \left( \sum_{kl} |\sigma_{ik}|^4 |G_{kl}^{(ij\mathbf{T})}|^2 \right)^{1/2} \left( \sum_{kl} |G_{kl}^{(ij\mathbf{T})}|^2 |\sigma_{lj}|^4 \right)^{1/2}. \quad (4.26)$$

Denote  $u_\alpha^{(ij\mathbf{T})}$  and  $\lambda_\alpha^{(ij\mathbf{T})}$  ( $\alpha = 1, 2, \dots, N - |\mathbf{T}| - 2$ ) the  $l^2$ -normalized eigenvectors and eigenvalues of  $H^{(ij\mathbf{T})}$ . Let  $u_\alpha^{(ij\mathbf{T})}(l)$  denote the  $l$ -th coordinate of  $u_\alpha^{(ij\mathbf{T})}$ , then for any  $l$

$$\sum_k |G_{kl}^{(ij\mathbf{T})}|^2 = \left( |G^{(ij\mathbf{T})}|^2 \right)_{ll} = \sum_\alpha \frac{|u_\alpha^{(ij\mathbf{T})}(l)|^2}{|\lambda_\alpha^{(ij\mathbf{T})} - z|^2} = \frac{\text{Im } G_l^{(ij\mathbf{T})}(z)}{\eta}. \quad (4.27)$$

Here we defined  $|A|^2 := A^* A$  for any matrix  $A$ . Inserting (4.27) into (4.25) and using the definition of  $M$  in (2.3), we obtain that (4.14) holds with a probability larger than  $1 - CN^{-c(\log \log N)}$  for sufficiently large  $N$ .

□

*Proof of Lemma 3.4.* We first prove (3.18) in the range  $3 \leq \eta \leq 10$ . Recall that  $\Omega^0$  is the subset of the entire probability space where  $\|H\| \leq 3$  see (3.11). By (7.11) from Lemma 7.2 (using that  $M \geq (\log N)^9$  is implied by (3.3)), we have  $\mathbb{P}(\Omega^0) \geq 1 - N^{-c \log \log N}$ . Denote  $\lambda_\alpha$  and  $\mathbf{u}_\alpha$  the eigenvalues and eigenvectors of  $H = (h_{ij})$ . From the identity

$$G_{ii} = \sum_\alpha \frac{|u_\alpha(i)|^2}{\lambda_\alpha - z} \quad (4.28)$$

and  $\max_\alpha |\lambda_\alpha| \leq 3$ , we have that

$$\eta^{-1} \geq |G_{ii}| \geq |\operatorname{Im} G_{ii}| \geq \frac{\eta}{(|E| + 3)^2 + \eta^2} \quad (4.29)$$

holds in  $\Omega^0$ . Together with  $3 \leq \eta \leq 10$  and  $|E| \leq 10$ , we obtain

$$c \leq |G_{ii}| \leq C \quad (4.30)$$

with some positive constants. From the interlacing property of the eigenvalues of the matrix and its submatrices, we find that not only  $\|H\| \leq 3$  but also  $\|H^{(\mathbf{T})}\| \leq 3$  holds on the set  $\Omega^0$ . Thus for any  $j, k$  such that  $i, j$  and  $k$  are all different, the bounds

$$c \leq |G_{ii}|, \quad |G_{ii}^{(j)}|, \quad |G_{ii}^{(jk)}| \leq C. \quad (4.31)$$

hold in  $\Omega^0$  by a similar argument that led to (4.30). Thus (4.6) implies

$$\mathbf{1}(\Omega^0) |G_{ij}^{(ij)}| \leq C^2 \mathbf{1}(\Omega^0) |K_{ij}^{(ij)}| \quad (4.32)$$

and (3.18) follows make use of (4.14) and  $\eta > 3$ .

Now we prove (3.19). Recall that the self consistent equation (3.16) with the error term  $\Upsilon_i(z)$  is given by (4.12), i.e.,

$$\Upsilon_i(z) = \sigma_{ii}^2 G_{ii} + \sum_{j \neq i} \sigma_{ij}^2 G_{ij} G_{ji} G_{ii}^{-1} + \left( K_{ii}^{(i)} - \mathbb{E}_{\mathbf{a}^i} K_{ii}^{(i)} \right). \quad (4.33)$$

Now we bound  $\Upsilon_i(z)$  in  $\widehat{\Omega}_z^o$ . Since  $\sigma_{ii}^2 \leq M^{-1}$ , with (4.31), the first term of the r.h.s. of (4.33) is less than  $O(M^{-1})$ . Then with (2.1), and using the bound on  $G_{ij}$  ( $i \neq j$ ) from (3.13) and the one on  $G_{ii}$  from (4.31), we obtain that the second term of the r.h.s. of (4.33) is less than  $C(\log N)^{10+8\alpha} (M\eta)^{-1}$  (and with (3.31), we know it is much less than 1), i.e., in  $\widehat{\Omega}_z^o$

$$\left| \sigma_{ii}^2 G_{ii} + \sum_{j \neq i} \sigma_{ij}^2 G_{ij} G_{ji} G_{ii}^{-1} \right| \leq C(\log N)^{10+8\alpha} (M\eta)^{-1}. \quad (4.34)$$

The last term of the r.h.s. of (4.33) can be bounded, using (4.13) with  $\mathbf{T} = \emptyset$ , with a very large probability. Using (4.8) and (4.31), the  $G_{kl}^{(i)}$ 's in (4.13) can be bounded as

$$|G_{kl}^{(i)}| \leq |G_{kl}| + C|G_{ki}||G_{il}|. \quad (4.35)$$

Therefore, again with the bound on  $G_{ij}$  ( $i \neq j$ ) in (3.13) and the one on  $G_{ii}$  from (4.31), we see that

$$|K_{ii}^{(i)} - \mathbb{E}_{\mathbf{a}^i} K_{ii}^{(i)}| \leq 2(\log N)^{8+6\alpha} (M\eta)^{-1/2} \quad (4.36)$$

holds in  $\widehat{\Omega}_z^o$  with a probability larger than  $\mathbb{P}(\widehat{\Omega}_z^o) - CN^{-c(\log \log N)}$  for sufficiently large  $N$ . Inserting (4.34) and (4.36) into (4.33) and together  $\eta \geq 3$ , we have proved (3.19).

Now we prove (3.20) for  $\eta \geq 3$ . By the definition of  $\Upsilon_i$  in (4.33), we have

$$\left| \mathbb{E} \left[ \mathbf{1}(\widehat{\Omega}_z^o) \Upsilon_i(z) \right] \right| \leq \mathbb{E} \mathbf{1}(\widehat{\Omega}_z^o) \left| \sigma_{ii}^2 G_{ii} + \sum_{j \neq i} \sigma_{ij}^2 G_{ij} G_{ji} G_{ii}^{-1} \right| + \left| \mathbb{E} \mathbf{1}([\widehat{\Omega}_z^o]^c) \left( K_{ii}^{(i)} - \mathbb{E}_{\mathbf{a}^i} K_{ii}^{(i)} \right) \right|, \quad (4.37)$$

since  $\mathbb{E} \left( K_{ii}^{(i)} - \mathbb{E}_{\mathbf{a}^i} K_{ii}^{(i)} \right) = 0$ . Using (4.31), in  $\Omega^0$  we have that  $|\sigma_{ii}^2 G_{ii} + \sum_{j \neq i} \sigma_{ij}^2 G_{ij} G_{ji} G_{ii}^{-1}|$  is always less than a constant  $C$  for some  $C > 0$ . Inserting this and (4.34) into (4.37), we obtain that

$$\left| \mathbb{E} \left[ \mathbf{1}(\widehat{\Omega}_z^o) \Upsilon_i(z) \right] \right| \leq C(\log N)^{10+8\alpha} (M\eta)^{-1} + \left| \mathbb{E} \left[ \mathbf{1}([\widehat{\Omega}_z^o]^c) |K_{ii}^{(i)} - \mathbb{E}_{\mathbf{a}^i} K_{ii}^{(i)}| \right] \right| + CN^{-c \log \log N}. \quad (4.38)$$

We now claim that for some large enough  $C > 0$  there exists  $c > 0$  such that

$$\mathbb{P}(|Z_{ii}^{(i)}| \geq N^C) \leq e^{-N^c} \quad \text{and} \quad \mathbb{P}(|K_{ii}^{(i)}| \geq N^C) \leq e^{-N^c}. \quad (4.39)$$

The first estimate follows from the definition of  $Z_{ii}^{(i)}$  given in Definition 4.1 by using the sub-exponential decay of the matrix elements and by using the trivial bound  $|G_{kl}^{(i)}| \leq \eta^{-1} \neq N$ . The second estimate is a trivial consequence of the first one and the definition of  $K_{ii}^{(i)}$ . Together with (4.38), we obtain (3.20) in the case that  $3 \leq \eta \leq 10$ .

We now prove (3.21) and (3.22) for the case  $\eta \leq 3$  satisfying (3.3). We will work in the event  $\Omega_{z'}^d \cap \Omega_z^d$ , where  $z' = z + iN^{-5}$ . Similarly as we proved (3.33), from the bound below (3.31) and the Lipschitz continuity of  $g(z)$ , we obtain that

$$|G_{ii}(z) - m_{sc}(z)| \leq 2(\log N)^{11+6\alpha} \frac{(\kappa + \eta)^{1/4}}{\sqrt{M\eta} g(z)} \quad (4.40)$$

and

$$|G_{ij}(z)| \leq 2(\log N)^{5+4\alpha} \frac{(\kappa + \eta)^{1/4}}{\sqrt{M\eta}} \quad (4.41)$$

hold in  $\Omega_{z'}^d \cap \Omega_z^d$ . We note the r.h.s of these inequalities are much less than  $(\log N)^{-1}$  by (3.3). From the explicit formula (2.13) we obtain that  $c \leq |m_{sc}(z)| \leq C$  for any  $|z| \leq 10$  with some positive constants. Using this observation and the fact that the r.h.s. of (4.40) is much less than  $(\log N)^{-1}$ , we have

$$c \leq |G_{ii}(z)| \leq C.$$

Hence, using (3.31), (4.7), (4.8) and the lower bound of  $|G_{ii}|$ , one can easily obtain that

$$|G_{ii}(z) - m_{sc}(z)|, \quad |G_{ii}^{(j)}(z) - m_{sc}(z)|, \quad |G_{ii}^{(jk)}(z) - m_{sc}(z)| \leq C(\log N)^{11+6\alpha} \frac{(\kappa + \eta)^{1/4}}{\sqrt{M\eta} g(z)} \quad (4.42)$$

hold in  $\Omega_{z'}^o \cap \Omega_z^d$ , (for the third term in l.h.s., we have also used the lower bounds of  $G_{ii}^{(j)}$ 's as above). Then we also have

$$c \leq |G_{ii}(z)|, \quad |G_{ii}^{(j)}(z)|, \quad |G_{ii}^{(jk)}(z)| \leq C \quad (4.43)$$

with some positive constants.

The definition of  $m_{sc}(z)$  implies  $\text{Im } m_{sc}(z) \leq C\sqrt{\kappa + \eta}$ . Then with (4.42), (3.3) and  $g(z) \leq \sqrt{\kappa + \eta}$ , we have that

$$\text{Im } G_{ii}^{(jk)}(z) \leq C\sqrt{\kappa + \eta} \quad (4.44)$$

holds in  $\Omega_z^o \cap \Omega_z^d$ , for some constant  $C > 0$ . Inserting it into (4.14), we obtain that

$$|K_{ij}^{(ij)}(z)| \leq C(\log N)^{4+4\alpha} \frac{(\kappa + \eta)^{1/4}}{\sqrt{M\eta}} \quad (4.45)$$

hold in  $\Omega_z^o \cap \Omega_z^d$ , with a probability larger than  $\mathbb{P}(\Omega_z^o \cap \Omega_z^d) - CN^{-c(\log \log N)}$  for sufficiently large  $N$ . Again, with (4.6) and (4.43), we obtain (3.21) for sufficiently large  $N$ .

Then, as we proved in (4.34) and (4.36), we get that

$$|\sigma_{ii}^2 G_{ii}(z) + \sum_{j \neq i} \sigma_{ij}^2 G_{ij} G_{ji} G_{ii}^{-1}(z)| \leq C(\log N)^{10+8\alpha} (M\eta)^{-1} \quad (4.46)$$

and

$$|K_{jj}^{(j)}(z) - \mathbb{E}_{\mathbf{a}^j} K_{jj}^{(j)}(z)| \leq C(\log N)^{8+6\alpha} (M\eta)^{-1/2} \quad (4.47)$$

hold in  $\widehat{\Omega}_z^o \cap \Omega_z^o \cap \Omega_z^d$ , with a probability larger than  $\mathbb{P}(\widehat{\Omega}_z^o \cap \Omega_z^o \cap \Omega_z^d) - CN^{-c(\log \log N)}$ , which implies (3.22).

Finally, similarly as using (4.37)- (4.38) to prove (3.20), we can obtain (3.23) in the case that  $\eta < 3$ .  $\square$

## 5 Stability of the self-consistent equation: proof of Lemma 3.5

In this section, we prove Lemma 3.5, i.e., we will prove the stability of the self-consistent equation with a precise error estimate given in (3.25). We set  $m_{sc} = m_{sc}(z)$  and  $\Upsilon = \max_i |\Upsilon_i(z)|$  for simplicity of notation and we will omit all  $z$  dependences in all the symbols. With the definition of  $m_{sc}(z)$  in (2.12) and (2.13), the following properties of  $m_{sc}(z)$  can be easily established:

**Lemma 5.1** *Let  $z = E + i\eta$  with  $\eta > 0$  and  $|z| \leq 20$ . Then we have*

$$|z + m_{sc}|^{-2} = |m_{sc}|^2 \leq 1 \quad (5.1)$$

and

$$|(z + m_{sc}(z))^{-2} - 1| \geq C\sqrt{\kappa + \eta} \quad (5.2)$$

for some constant  $C$ . Furthermore, suppose that either  $2 \leq |E| \leq 10$  or  $\kappa \leq \eta$ . Then

$$|z + m_{sc}|^{-2} = |m_{sc}|^2 \leq 1 - C\sqrt{\kappa + \eta}. \quad (5.3)$$

For small values of  $|z^2 - 4| \asymp \kappa + \eta$ ,  $m_{sc}(z)$  has the asymptotic expansion

$$m_{sc} = \mp 1 + \frac{1}{2}\sqrt{z^2 - 4} + O(|z^2 - 4|), \quad \text{near } z \asymp \pm 2. \quad (5.4)$$

$\square$

We first prove (3.25) for the case that  $3 \leq \eta \leq 10$ . In this case, we can easily check that  $g(z) = \sqrt{\kappa + \eta}$ . Denote the difference between  $G_{ii}$  and  $m_{sc}$  by

$$v_i = G_{ii} - m_{sc}, \quad 1 \leq i \leq N.$$

By the self consistent equation (3.16), (2.1) and (2.12), we have

$$v_i = \frac{\sum_i \sigma_{ij}^2 v_j + \Upsilon_i}{(z + m_{sc} + \sum_j \sigma_{ij}^2 v_j + \Upsilon_i)(z + m_{sc})}, \quad 1 \leq i \leq N. \quad (5.5)$$

For  $\eta \geq 3$ ,  $|z + m_{sc}(z)| > 2$  by (2.13). Using  $|G_{ii}| \leq \eta^{-1}$  and  $|m_{sc}| \leq \eta^{-1}$ , we obtain

$$|v_i| \leq 2/\eta \leq 2/3, \quad 1 \leq i \leq N. \quad (5.6)$$

From the assumption (3.17) and (3.3), we have  $\Upsilon = \max_i |\Upsilon_i| \ll 1$  in this region. Together with  $|z + m_{sc}(z)| > 2$  and (5.6), we obtain that the absolute value of the r.h.s. of (5.5) is less than

$$\frac{\sup_i |v_i|}{|z + m_{sc}(z)| - \sup_i |v_i|} + O(\Upsilon). \quad (5.7)$$

Taking the absolute value of (5.5) and maximizing over  $n$ , we have

$$\sup_n |v_n| \leq \frac{\sup_i |v_i|}{|z + m_{sc}| - \sup_i |v_i|} + O(\Upsilon). \quad (5.8)$$

The denominator satisfies  $|z + m_{sc}(z)| - \sup_i |v_i| \geq 2 - 2/3 = 4/3$ , therefore we obtain  $\sup |v_i| = \sup_i |G_{ii} - m_{sc}(z)| \leq O(\Upsilon)$ , which shows (3.25) for  $3 \leq \eta \leq 10$ .

Next, we prove (3.25) in the case that  $\eta \leq 3$  with  $\eta$  satisfying (3.3) and under the condition (3.24). Define

$$m = m(z) := \frac{1}{N} \sum_i G_{ii} \quad \text{and} \quad u_i := G_{ii} - m. \quad (5.9)$$

Combining (3.17), (3.3), (3.24) with the fact that  $g(z) \leq C$ , we can see that

$$\Upsilon \leq (\log N)^{-4} g^2(z) \leq C(\log N)^{-4}, \quad |G_{ii} - m_{sc}(z)| \leq C(\log N)^{-2}. \quad (5.10)$$

Together with (5.1), we have

$$|z + m_{sc}(z)| - |G_{ii} - m_{sc}(z)| - |\Upsilon| \geq C$$

for some  $C > 0$ . Furthermore (3.24) implies

$$|m(z) - m_{sc}| \leq 2(\log N)^{-2} g(z) \quad (5.11)$$

thus there exists  $C > 0$  such that

$$|z + m(z)| - |G_{ii} - m(z)| - |\Upsilon| \geq C. \quad (5.12)$$

Therefore, expanding the self consistent equation (3.16) around  $z + m(z)$ , we obtain that

$$0 = G_{ii} + \frac{1}{z + \sum_j \sigma_{ij}^2 G_{jj} + \Upsilon_i} = G_{ii} + \frac{1}{z + m(z)} + \Omega_i \quad (5.13)$$

where  $\Omega_i$  is defined by the second equality and it satisfies

$$\Omega_i = -\frac{\sum_j \sigma_{ij}^2 u_j}{(z + m(z))^2} + O(\|\mathbf{u}\|_\infty^2) + O(\Upsilon) \quad (5.14)$$

with error bounds uniform in  $i$ . Here  $\|\mathbf{u}\|_\infty = \max_i |u_i|$ . Taking the average of the r.h.s of (5.13) with respect to  $i$ , we obtain that

$$m(z) + \frac{1}{z + m(z)} = -\Omega \quad (5.15)$$

where

$$\Omega := \frac{1}{N} \sum_i \Omega_i, \quad (5.16)$$

and it satisfies

$$|\Omega| \leq O(\|\mathbf{u}\|_\infty^2) + O(\Upsilon). \quad (5.17)$$

Here we used  $\sum_i \sum_j \sigma_{ij}^2 u_j = \sum_j u_j = 0$ . The bound (3.24), (5.11) and  $g(z) \leq \sqrt{\kappa + \eta}$  (from (3.2)) implies that

$$\|\mathbf{u}\|_\infty \leq 4(\log N)^{-2} g(z) \leq C(\log N)^{-2} \sqrt{\kappa + \eta}. \quad (5.18)$$

Together with (5.10) and (5.17), we obtain

$$|\Omega| \leq C(\log N)^{-4} (\kappa + \eta).$$

To bound  $m(z)$ , we use the following lemma.

**Lemma 5.2** *Let  $z = E + i\eta \in \mathbb{C}$ ,  $|z| \leq 10$  and let  $\delta > 0$  be a sufficiently small constant. Let  $t \in \mathbb{C}$  such that*

$$|t| \leq \delta(\kappa + \eta). \quad (5.19)$$

*Suppose there is a function  $s_z(t) \in \mathbb{C}$  that solves the equation*

$$s_z(t) + \frac{1}{z + s_z(t)} = t, \quad (5.20)$$

*with  $\text{Im} s_z(t) > 0$  and the estimate*

$$|s_z(t) - m_{sc}(z)| \leq \delta \sqrt{\kappa + \eta} \quad (5.21)$$

*holds. Then*

$$|s_z(t) - m_{sc}(z)| \leq C \frac{|t|}{\sqrt{\kappa + \eta}} \quad (5.22)$$

*for some constant  $C > 0$ .*

*Proof of Lemma 5.2.* It follows from (5.20) that

$$s_z(t) = t + \frac{-z - t}{2} \pm \frac{\sqrt{(z + t)^2 - 4}}{2}. \quad (5.23)$$

We denote by  $s_z^1(t)$  and  $s_z^2(t)$  the two solutions of this equation, which are continuous with respect to  $t$  locally in the neighborhood (5.19). When  $t = 0$ , one of them is equal to  $m_{sc}(z)$ , we choose  $s_z^1(0) = m_{sc}(z)$ . From (5.23), we have

$$|s^1 - s^2| = |(z + t)^2 - 4|^{1/2}. \quad (5.24)$$

Then, for small enough  $\delta$ , if  $|t| \leq \delta(\kappa + \eta)$ , then  $|s^1 - s^2| \geq \frac{1}{2} \min\{|z - 2|, |z + 2|\}$  by using (5.24) and that  $\kappa + \eta \asymp \min\{|z - 2|, |z + 2|\}$ . We thus see that only one out of  $s^1$  and  $s^2$  can satisfy (5.21). With the assumption that  $s_z^1(0) = m_{sc}(z)$ , it is  $s^1$  that satisfies (5.21). Then

$$s_z(t) - m_{sc}(z) = s_z^1(t) - s_z^1(0) \quad (5.25)$$

and (5.22) follows from the fact that

$$|\partial_t s_z(t)| \leq O\left(\frac{1}{|\sqrt{(z+t)^2 - 4}|}\right) \leq O\left(\frac{1}{\sqrt{\kappa + \eta}}\right), \quad (5.26)$$

where for the second inequality, we used  $|t| \leq \delta(\kappa + \eta)$ .  $\square$

Using Lemma 5.2, for  $s_z(t) = m(z)$  and  $t = -\Omega$ , we have

$$|m(z) - m_{sc}(z)| \leq \frac{C|\Omega|}{\sqrt{\kappa + \eta}} \leq C(\log N)^{-2} \|\mathbf{u}\|_\infty + \frac{C\Upsilon}{\sqrt{\kappa + \eta}}, \quad (5.27)$$

where in the second inequality we used (5.17) and (5.18). Subtracting (5.17) from (5.13), we have the equation for  $u_i$

$$u_i = G_{ii} - m(z) = \frac{\sum_j \sigma_{ij}^2 u_j}{(z + m(z))^2} + \Omega + O(\|\mathbf{u}\|_\infty^2) + O(\Upsilon) = w_i + \frac{\sum_j \sigma_{ij}^2 u_j}{(z + m_{sc}(z))^2}, \quad (5.28)$$

where  $w_i$  is defined as  $u_i - (\sum_j \sigma_{ij}^2 u_j)(z + m_{sc})^{-2}$ . By (5.17), it is bounded by

$$\|\mathbf{w}\|_\infty = O(\|\mathbf{u}\|_\infty^2) + O(|\|\mathbf{u}\|_\infty|(z + m)^{-2} - (z + m_{sc})^{-2}|) + O(\Upsilon). \quad (5.29)$$

Then, using (5.11) and (5.1), we obtain that

$$|(z + m)^{-2} - (z + m_{sc})^{-2}| \leq C|m(z) - m_{sc}(z)|. \quad (5.30)$$

Inserting this into (5.29), using the bounds on  $\|\mathbf{u}\|_\infty$  in (5.18) and (5.27), we have

$$\|\mathbf{w}\|_\infty = O(\|\mathbf{u}\|_\infty^2) + O(\Upsilon). \quad (5.31)$$

From (5.3) in Lemma 5.1, whenever  $|E| \geq 2$  or  $\kappa \leq \eta$ , in which case  $g(z) \asymp \sqrt{\kappa + \eta}$ , we have

$$|z + m_{sc}|^{-2} \leq 1 - C\sqrt{\kappa + \eta}, \quad (5.32)$$

for some  $C > 0$ . Therefore (5.28) imply in this region that

$$\|\mathbf{u}\|_\infty \leq C(\kappa + \eta)^{-1/2} \|\mathbf{w}\|_\infty. \quad (5.33)$$

Using (5.31), we get

$$\|\mathbf{u}\|_\infty \leq \frac{C}{\sqrt{\kappa + \eta}} \|\mathbf{u}\|_\infty^2 + \frac{C}{\sqrt{\kappa + \eta}} \Upsilon, \quad (5.34)$$

and using  $\|\mathbf{u}\|_\infty \ll \sqrt{\kappa + \eta}$  from (5.18), we conclude that

$$\|\mathbf{u}\|_\infty \leq O\left(\frac{\Upsilon}{\sqrt{\kappa + \eta}}\right). \quad (5.35)$$



Combining this with the bound on  $m - m_{sc}$  (5.27) and  $G_{ii} - m_{sc} = m - m_{sc} + u_i$ , we obtain (3.25).

Finally, we consider the main interesting regime:  $|E| \leq 2$  and  $\kappa \geq \eta$ . We claim that the following inequality about  $m_{sc}(z)$  holds.

**Lemma 5.3** *Let  $1 > \delta_- > 0$  be a given constant. Then there exist small real numbers  $\tau \geq 0$  and  $c_1 > 0$ , depending only on  $\delta_-$ , such that we have*

$$\max_{x \in [-1+\delta_-, 1-\delta_+]} \left\{ \left| \tau + x m_{sc}^2 \right|^2 \right\} \leq (1 - c_1 \widehat{g}(z)) (1 + \tau)^2 \quad (5.36)$$

with

$$\widehat{g}(z) = \max\{\delta_+, |1 - \operatorname{Re} m_{sc}^2(z)|\} \quad (5.37)$$

for any positive number  $\delta_+$  such that  $-1 + \delta_- \leq 1 - \delta_+$ .

We postpone the proof of this lemma to the end of this subsection and we first complete the main argument. Recall that  $B = \{\sigma_{ij}^2\}_{i,j=1}^N$  is the matrix of variances which is symmetric. We also recall  $\delta_{\pm}$  from (2.4) and we will apply Lemma 5.3 with these  $\delta_-$  and  $\delta_+$ . Fix  $z$ , set  $\zeta := m_{sc}^2(z) = (m_{sc}(z) + z)^{-2}$  and rewrite (5.28) as

$$\mathbf{u} = (I - \zeta B)^{-1} \mathbf{w} = \frac{1}{1 + \tau} \left[ I - \frac{\zeta B + \tau}{1 + \tau} \right]^{-1} \mathbf{w} \quad (5.38)$$

with  $\tau$  given in Lemma 5.3. Define  $Q := I - |\mathbf{e}\rangle\langle\mathbf{e}|$  to be the projection onto the orthogonal complement of the normalized eigenvector  $\mathbf{e} = N^{-1/2}(1, 1, \dots, 1)$  belonging to the simple eigenvalue 1 of  $B$ . Note that  $B$  and  $Q$  commute and that the spectrum of  $BQ$  lies in  $[-1 + \delta_-, 1 - \delta_+]$ . Denote by  $\|A\|$  the usual  $\ell^2 \rightarrow \ell^2$  norm of a matrix  $A$ . Since

$$\left\| \frac{\zeta B + \tau}{1 + \tau} Q \right\| \leq \sup_{x \in [-1+\delta_-, 1-\delta_+]} \left| \frac{\zeta x + \tau}{1 + \tau} \right| \leq (1 - c_1 \widehat{g}(z))^{1/2} < 1$$

by the Lemma 5.3 and  $\mathbf{w} \perp (1, \dots, 1)$ , the Neumann expansion of (5.38) converges on  $\operatorname{span}((1, \dots, 1)^\perp)$  and

$$\mathbf{u} = (I - \zeta B)^{-1} \mathbf{w} = \frac{1}{1 + \tau} \sum_{n=0}^{\infty} \left( \frac{\zeta B + \tau}{1 + \tau} \right)^n \mathbf{w}. \quad (5.39)$$

We will compute the  $\ell^\infty \rightarrow \ell^\infty$  norm of this matrix. First note that

$$\left\| \frac{\zeta B + \tau}{1 + \tau} \right\|_{\infty \rightarrow \infty} = \max_i \sum_j \left| \left( \frac{\zeta B + \tau}{1 + \tau} \right)_{ij} \right| \leq \frac{1}{1 + \tau} \max_i \sum_j |\zeta B_{ij} + \tau \delta_{ij}| \leq \frac{|\zeta| + \tau}{1 + \tau} \leq 1, \quad (5.40)$$

since  $|\zeta| = |m_{sc}|^2 \leq 1$  and  $\sum_j |B_{ij}| = \sum_j B_{ij} = \sum_j \sigma_{ij}^2 = 1$ . Then we have

$$\left\| \left( \frac{\zeta B + \tau}{1 + \tau} \right)^n \mathbf{u} \right\| = \left\| \left( \frac{\zeta B + \tau}{1 + \tau} \right)^n Q \mathbf{u} \right\| \leq \sup_{x \in [-1+\delta_-, 1-\delta_+]} \left| \frac{\zeta x + \tau}{1 + \tau} \right|^n \|\mathbf{u}\| \leq (1 - c_1 \widehat{g}(z))^{n/2} \|\mathbf{u}\|$$

by Lemma 5.3. Since for any  $N \times N$  matrix we have

$$\|A\|_{\infty \rightarrow \infty} \leq \sqrt{N} \|A\|,$$

we obtain

$$\left\| \left( \frac{\zeta B + \tau}{1 + \tau} \right)^n Q \right\|_{\infty \rightarrow \infty} \leq \sqrt{N} (1 - c_1 \widehat{g}(z))^{n/2}. \quad (5.41)$$

Thus, estimating the first  $n \leq n_0 := (\log N)(c_1 \widehat{g}(z))^{-1}$  terms in (5.39) by (5.40), and the rest by (5.41), we get

$$\|\mathbf{u}\|_\infty \leq \left( \frac{\log N}{c_1 \widehat{g}(z)} + \sum_{n=n_0}^{\infty} \sqrt{N} (1 - c_1 \widehat{g}(z))^{n/2} \right) \|\mathbf{w}\|_\infty \leq C \frac{\log N}{\widehat{g}(z)} \|\mathbf{w}\|_\infty.$$

Using the bound (5.31) on  $\|\mathbf{w}\|_\infty$  and the bound (5.18) on  $\|\mathbf{u}\|_\infty$ , we have

$$\|\mathbf{u}\|_\infty \leq (\log N)^{-1} \|\mathbf{u}\|_\infty + C \frac{\log N}{\widehat{g}(z)} \Upsilon \quad (5.42)$$

which implies

$$\|\mathbf{u}\|_\infty \leq C \frac{\log N}{\widehat{g}(z)} \Upsilon \quad (5.43)$$

for some  $C > 0$ . Combining this with (5.27), we find

$$|G_{ii} - m_{sc}| \leq |m(z) - m_{sc}(z)| + |u_i| \leq C \left[ \|\mathbf{u}\|_\infty + \frac{\Upsilon}{\sqrt{\kappa + \eta}} \right] \leq C \left[ \frac{\log N}{\widehat{g}(z)} + \frac{1}{\sqrt{\kappa + \eta}} \right] \Upsilon$$

which implies (3.25), since  $g(z) = \min\{\sqrt{\kappa + \eta}, \widehat{g}(z)\}$ .  $\square$

*Proof of Lemma 5.3.* First, if  $\widehat{g}(z) = \delta_+$ , then we choose  $\tau = 0$ . With  $|m_{sc}| \leq 1$  in (5.1), one can see that (5.36) holds.

In the case of  $\widehat{g}(z) = |1 - \operatorname{Re} m_{sc}^2|$ , we have  $\widehat{g} \leq 2$  by using  $|m_{sc}| \leq 1$ . We choose  $\tau = \delta_-/10$ , then

$$\max_{x \in [-1 + \delta_-, 1 - \delta_+]} \left| \tau + x m_{sc}^2 \right|^2 \leq \max \left\{ \left| \frac{\delta_-}{10} + m_{sc}^2 \right|^2, \left| \frac{\delta_-}{10} - (1 - \delta_-) m_{sc}^2 \right|^2 \right\}, \quad (5.44)$$

and

$$\left| \frac{\delta_-}{10} - (1 - \delta_-) m_{sc}^2 \right|^2 \leq \left| 1 - \frac{9\delta_-}{10} \right|^2 \leq (1 - \tau \widehat{g}(z))(1 + \delta_-/10)^2. \quad (5.45)$$

For the other term in r.h.s. of (5.44), we have

$$\left| \frac{\delta_-}{10} + m_{sc}^2 \right|^2 = |m_{sc}|^4 + 0.2\delta_- \operatorname{Re}(m_{sc}^2) + (\delta_-/10)^2. \quad (5.46)$$

With  $|m_{sc}| \leq 1$  in (5.1) and  $\widehat{g}(z) = |1 - \operatorname{Re}(m_{sc}^2)|$  in this case, (5.46) is bounded as

$$\left| \frac{\delta_-}{10} + m_{sc}^2 \right|^2 \leq \left( 1 + \frac{\delta_-}{10} \right)^2 - 0.2\delta_- \widehat{g}(z) \leq \left( 1 + \frac{\delta_-}{10} \right)^2 (1 - C\widehat{g}(z)) \quad (5.47)$$

for some  $C$  depending on  $\delta_-$ . At last, we complete the proof by combining (5.47) and (5.45).  $\square$

## 6 Proof of the universality of local statistics

We now outline the main steps to prove Theorem 2.2.

*Step 1. Local relaxation flow.* Following [19], we first prove that the local eigenvalue statistics of Dyson Brownian motion (DBM) at a fixed time  $t$  are the same as those of GUE if  $t \asymp N^{-\varepsilon_0}$  for some  $\varepsilon_0$ . The DBM is generated by the flow

$$H_t = e^{-t/2} H_0 + (1 - e^{-t})^{1/2} V, \quad (6.1)$$

where  $H_0$  is the initial matrix and  $V$  is an independent GUE matrix whose matrix elements are centered Gaussian random variables with variance  $1/N$ . Strictly speaking, for each matrix element we have used the Ornstein-Uhlenbeck (OU) process on  $\mathbb{C}$  instead of the Brownian motion which was used in the original definition of DBM in [19]. It is easy to check that the eigenvalues of  $H_t$  follow a process, very similar to the original DBM in [19], but with a drift. With a slight abuse of terminology, we will still call this process DBM. More precisely, let

$$\mu = \mu_N(d\mathbf{x}) = \frac{e^{-\mathcal{H}(\mathbf{x})}}{Z_\beta} d\mathbf{x}, \quad \mathcal{H}(\mathbf{x}) = N \left[ \beta \sum_{i=1}^N \frac{x_i^2}{4} - \frac{\beta}{N} \sum_{i < j} \log |x_j - x_i| \right] \quad (6.2)$$

( $\beta = 2$  for GUE) be the probability measure of the eigenvalues of the general  $\beta$  ensemble,  $\beta \geq 1$  (in this section, we often use the notation  $x_j$  for the eigenvalues to follow the notations of [19]). In this paper we consider the  $\beta = 2$  case for simplicity, but we stress that our proof applies to the case of symmetric matrices as well. Denote the distribution of the eigenvalues at time  $t$  by  $f_t(\mathbf{x})\mu(d\mathbf{x})$ . Then  $f_t$  satisfies

$$\partial_t f_t = \mathcal{L} f_t. \quad (6.3)$$

where (see (2.2) in [19])

$$\mathcal{L} = \sum_{i=1}^N \frac{1}{2N} \partial_i^2 + \sum_{i=1}^N \left( -\frac{\beta}{4} x_i + \frac{\beta}{2N} \sum_{j \neq i} \frac{1}{x_i - x_j} \right) \partial_i. \quad (6.4)$$

**Theorem 6.1** *Suppose that the probability law for the initial matrix  $H_0$  satisfies the assumptions of Theorem 2.2. Then there exists  $\varepsilon_0 > 0$  such that for any*

$$t \geq N^{-\varepsilon_0}, \quad (6.5)$$

*the probability law for the eigenvalues of  $H_t$  satisfies (2.20), i.e., for any  $k \geq 1$  and for any compactly supported continuous test function  $O : \mathbb{R}^k \rightarrow \mathbb{R}$ , we have*

$$\begin{aligned} & \lim_{\kappa \rightarrow 0} \lim_{N \rightarrow \infty} \frac{1}{2\kappa} \int_{E-\kappa}^{E+\kappa} dE' \int_{\mathbb{R}^k} d\alpha_1 \dots d\alpha_k O(\alpha_1, \dots, \alpha_k) \\ & \quad \times \frac{1}{\varrho_{sc}(E)^k} \left( p_N^{(k)} - p_{GUE,N}^{(k)} \right) \left( E' + \frac{\alpha_1}{N \varrho_{sc}(E)}, \dots, E' + \frac{\alpha_k}{N \varrho_{sc}(E)} \right) = 0, \end{aligned}$$

where  $p_{GUE,N}^{(k)}$  is the  $k$ -point correlation function of the GUE ensemble.

*Proof of Theorem 6.1.* We first recall the following general theorem concerning the Dyson Brownian motion from [19] that asserts that under four general assumptions, the local eigenvalue statistics of the time evolved matrix  $H_t$  coincide with GUE. The first assumption (called Assumption I in [19]) is a convexity bound on  $\mathcal{H}$  which is automatically satisfied in our case and we only have to verify the following three assumptions.

**Assumption II.** There exists a continuous, compactly supported density function  $\varrho(x) \geq 0$ ,  $\int_{\mathbb{R}} \varrho = 1$ , on the real line, independent of  $N$ , such that for any fixed  $a, b \in \mathbb{R}$

$$\lim_{N \rightarrow \infty} \sup_{t \geq 0} \left| \int \frac{1}{N} \sum_{j=1}^N \mathbf{1}(x_j \in [a, b]) f_t(\mathbf{x}) d\mu(\mathbf{x}) - \int_a^b \varrho(x) dx \right| = 0. \quad (6.6)$$

For the next assumption, we introduce a notation. Let  $\gamma_j = \gamma_{j,N}$  denote the location of the  $j$ -th point under the limiting density, i.e.,  $\gamma_j$  is defined by

$$N \int_{-\infty}^{\gamma_j} \varrho(x) dx = j, \quad 1 \leq j \leq N, \quad \gamma_j \in \text{supp} \varrho. \quad (6.7)$$

We will call  $\gamma_j$  the *classical location* of the  $j$ -th point.

**Assumption III.** There exists an  $\varepsilon > 0$  such that

$$\sup_{t \geq 0} \int \frac{1}{N} \sum_{j=1}^N (x_j - \gamma_j)^2 f_t(\mathbf{x}) \mu(d\mathbf{x}) \leq CN^{-1-2\varepsilon} \quad (6.8)$$

with a constant  $C$  uniformly in  $N$ .

The final assumption is an upper bound on the local density. For any  $I \in \mathbb{R}$ , let

$$\mathcal{N}_I := \sum_{i=1}^N \mathbf{1}(x_i \in I)$$

denote the number of points in  $I$ .

**Assumption IV.** For any compact subinterval  $I_0 \subset \{E : \varrho(E) > 0\}$  independent of  $N$ , and for any  $\delta > 0$ ,  $\sigma > 0$  and  $r > 0$ , there are constants  $c$  depending on  $I_0$ ,  $\delta$ ,  $\sigma$  and  $r$  such that for any interval  $I \subset I_0$  with  $|I| \geq N^{-1+\sigma}$ , we have

$$\sup_{\tau \geq N^{-2\varepsilon+\delta}} \int \mathbf{1}\{\mathcal{N}_I \geq KN|I|\} f_\tau d\mu \leq N^{-c \log \log N}, \quad K = N^r \quad (6.9)$$

where  $\varepsilon$  is the exponent from Assumption III.

**Theorem 6.2** [19, Theorem 2.1] *Let  $\varepsilon > 0$  be the exponent from Assumption III. Suppose that there is a time  $\tau < N^{-2\varepsilon}$  such that the following entropy bound holds*

$$S_\mu(f_\tau) := \int f_\tau \log f_\tau d\mu \leq CN^m \quad (6.10)$$

for some fixed exponent  $m$ . Suppose that the Assumptions II, III and IV hold for the solution  $f_t$  of the forward equation (6.3) for all time  $t \geq \tau$ . Let  $E \in \mathbb{R}$  be a point where  $\varrho(E) > 0$ . Then for any  $k \geq 1$  and for any compactly supported continuous test function  $O : \mathbb{R}^k \rightarrow \mathbb{R}$ , we have

$$\begin{aligned} \lim_{b \rightarrow 0} \lim_{N \rightarrow \infty} \sup_{t \geq N^{-2\varepsilon + \delta}} \frac{1}{2b} \int_{E-b}^{E+b} dE' \int_{\mathbb{R}^k} d\alpha_1 \dots d\alpha_k O(\alpha_1, \dots, \alpha_k) \\ \times \frac{1}{\varrho(E)^k} \left( p_{t,N}^{(k)} - p_{\mu,N}^{(k)} \right) \left( E' + \frac{\alpha_1}{N\varrho(E)}, \dots, E' + \frac{\alpha_k}{N\varrho(E)} \right) = 0. \end{aligned} \quad (6.11)$$

Theorem 6.2 was exactly Theorem 2.1 of [19] except that the assumption (6.10) on the entropy in [19] was stated for the initial probability density  $f_0$ . Clearly, we can start the flow (6.4) from a fixed time  $\tau \ll N^{-2\varepsilon + \delta}$  since the statement of Theorem 6.2 concerns only the time  $t \geq N^{-2\varepsilon + \delta}$ . In the case that the flow (6.4) is generated from the matrix evolution (6.1), the entropy assumption (6.10) is satisfied automatically. To see this, let  $\nu_t^{ij}$  denote the probability measure of the  $ij$ -th element of the matrix  $H_t$ ,  $i \leq j$ , and  $\bar{\nu}_t$  the probability measure of the matrix  $H_t$ . Let  $\bar{\mu}$  denote the probability measure of the GUE and  $\mu^{ij}$  the probability measure of its  $ij$ -th element which is a Gaussian measure with mean zero and variance  $1/N$ . Since the dynamics of matrix elements are independent (subject to the Hermitian condition), we have the identity

$$\int \log \left( \frac{d\bar{\nu}_t}{d\bar{\mu}} \right) d\bar{\nu}_t = \sum_{ij} \int \log \left( \frac{d\nu_t^{ij}}{d\mu^{ij}} \right) d\nu_t^{ij}. \quad (6.12)$$

The process  $t \rightarrow \nu_t^{ij}$  is an Ornstein-Uhlenbeck process and each entropy term on the right hand side of the last equation is bounded by  $CN$  provided that  $t \geq 1/N$  and  $\nu_0^{ij}$  has a subexponential decay. It is easy to check from the explicit OU kernel. Since the entropy of the marginal distribution on the eigenvalues is bounded by the entropy of the total measure on the matrix, we have proved that

$$\int f_{1/N} \log f_{1/N} d\mu \leq CN^3, \quad (6.13)$$

and this verifies (6.10). Therefore, in order to apply Theorem 6.2, we only have to verify the Assumptions II, III and IV. Clearly, Assumption II follows from Theorem 2.1 (note that in the case of generalized Wigner matrix,  $M \asymp N$  and  $g(z) \asymp \sqrt{\kappa + \eta}$ ). Assumption IV also follows from Theorem 2.1 by noting that  $N_I \leq C \text{Im}(E + i\eta)$  if  $I$  is an interval of length  $\eta$  about  $E$ . We also note that Assumption IV in [19] was stated in a slightly stronger form, requiring a large deviation bound (6.9) for all  $K \geq 1$ , but inspecting the proof of Theorem 2.1 of [19] reveals that Assumption IV is used only for  $K$  larger than some positive power of  $N$  and smaller than  $N$  (the main observation is that the upper limit of the summation in (7.16) of [19] is effectively  $N$  and not  $\infty$ ).

Having verified all other assumptions, it remains to prove (6.8), which we state as the next theorem.

**Theorem 6.3** *Suppose  $H$  satisfies the assumptions of Theorem 2.2, in particular, it is a generalized Wigner matrix with positive constants  $C_{inf}$ ,  $C_{sup}$  in (2.6). Let  $\tilde{\nu}_{ij}(x)dx := \sigma_{ij}\nu_{ij}(\sigma_{ij}x)dx$  be the rescaling of the distributions  $\nu_{ij}$  of the matrix elements and suppose that they satisfy the logarithmic Sobolev inequality (LSI) with a constant  $C_S$  independent of  $N, i, j$ , i.e.,*

$$\int u \log u d\tilde{\nu}_{ij} \leq C_S \int |\nabla \sqrt{u}|^2 d\tilde{\nu}_{ij} \quad (6.14)$$

holds for any smooth probability density  $u$ ,  $\int u d\tilde{\nu}_{ij} = 1$ . Denote  $\lambda_i$  the  $i$ -th eigenvalue of  $H$  in increasing order,  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N$ . Then there exists  $\varepsilon > 0$  depending on  $\alpha, \beta$  in (2.14) but independent of  $C_{inf}$ ,  $C_{sup}$  and  $C_S$  such that

$$\frac{1}{N} \sum_{i=1}^N \mathbb{E}(\lambda_i - \gamma_i)^2 \leq CN^{-1-2\varepsilon}, \quad (6.15)$$

if  $N$  is sufficiently large (depending on  $C_{inf}$ ,  $C_{sup}$ ,  $C_S$ ,  $\alpha$  and  $\beta$ ).

The proof of Theorem 6.3 will be given in Section 7. It is easy to check that if an initial matrix  $H = H_0$  satisfies the conditions of Theorem 6.3, then its evolution  $H_t$  under the Ornstein-Uhlenbeck flow will also satisfy these conditions with constants changed at most by a factor two. The main condition to check is that the logarithmic Sobolev inequality (6.14) holds for  $0 \leq t \ll 1$ . But this was proved in the argument following Lemma 5.3 of [15] using an estimate on the logarithmic Sobolev constant for convolution of two measures, i.e., Lemma B.1 of [15]. Therefore Theorem 6.3 guarantees (6.15) for all positive times  $t > 0$  and this proves Assumption III provided that the initial distribution satisfies the LSI (6.14). We have thus proved Theorem 2.2 for matrix ensembles of the form

$$h_{ij} = e^{-t/2} \widehat{h}_{ij} + (1 - e^{-t})^{1/2} N^{-1/2} \xi_{ij}^G, \quad t \geq N^{-2\varepsilon + \delta} \quad (6.16)$$

where  $\xi_{ij}^G$  are i.i.d. complex random variables with Gaussian distribution with mean 0 and variance 1, and  $\widehat{h}_{ij}$ 's are independent random variables such that the rescaled variables  $\widehat{\zeta}_{ij} = \widehat{h}_{ij}/\sigma_{ij}$  satisfy the LSI assumption (6.14). In (6.16)  $\delta > 0$  is arbitrary and  $\varepsilon$  is fixed in Theorem 6.3. In particular, with the choice  $\delta = \varepsilon$  and  $t \asymp N^{-\varepsilon}$ , we have proved Theorem 2.2 for matrix ensembles  $h_{ij} = \sigma_{ij} \zeta_{ij}$  if  $\zeta_{ij}$  is of the form

$$\zeta_{ij} = (1 - \gamma)^{1/2} \widehat{\zeta}_{ij} + \gamma^{1/2} \xi_{ij}^G, \quad \gamma \asymp N^{-\varepsilon}, \quad \text{distribution of } \widehat{\zeta}_{ij} \text{ satisfies (6.14)}. \quad (6.17)$$

*Step 2. Eigenvalue correlation function comparison theorem.*

The next step is to prove that the correlation functions of eigenvalues for two matrix ensembles are identical up to scale  $1/N$  provided that the first four moments of all matrix elements of these two ensembles are almost identical. This theorem is a corollary of Theorem 2.3 and we state it as the following correlation function comparison theorem. The proof will be given in Section 8. Note that the assumption (2.21) in Theorem 2.3 is satisfied by Theorem 3.1; in case of generalized Wigner matrix we have  $g(z) = \sqrt{\kappa + \eta}$  and  $M \asymp N$ , so in the regime where  $|E|$  is separated away from 2, we have from (3.4), that  $G_{ii}(z)$  is uniformly bounded (modulo logarithmic factors).

**Theorem 6.4** *Suppose the assumptions of Theorem 2.3 hold. Let  $p_{v,N}^{(k)}$  and  $p_{w,N}^{(k)}$  be the  $k$ -point functions of the eigenvalues w.r.t. the probability law of the matrix  $H^{(v)}$  and  $H^{(w)}$ , respectively. Then for any  $|E| < 2$ , any  $k \geq 1$  and any compactly supported continuous test function  $O : \mathbb{R}^k \rightarrow \mathbb{R}$  we have*

$$\int_{\mathbb{R}^k} d\alpha_1 \dots d\alpha_k O(\alpha_1, \dots, \alpha_k) \left( p_{v,N}^{(k)} - p_{w,N}^{(k)} \right) \left( E + \frac{\alpha_1}{N}, \dots, E + \frac{\alpha_k}{N} \right) = 0. \quad (6.18)$$

*Step 3. Approximation of a measure by Ornstein-Uhlenbeck process for small time.*

Summarizing, we have proved Theorem 2.2 in Step 1 for matrix ensembles whose probability distributions of the normalized matrix elements  $\zeta_{ij}$  are of the form (6.17). Using the Green's function comparison theorem,

i.e. Theorem 6.4, we extended the class of distributions to all random variables whose first four moments can almost be matched (more precisely, match the first three moments and almost match the fourth moments in the sense of (2.22)) by random variables in the class (6.17). In order to complete the proof of Theorem 2.2, it remains to prove that for all measures in the class given by the assumptions of Theorem 2.2, i.e., measures satisfying the subexponential decay condition, the uniformly bounded-variance condition (2.6) and the moment restriction (2.19) for the real and imaginary parts, we can find random variables in the class (6.17) to almost match the first four moments. Since the real and imaginary parts are i.i.d., it is sufficient to match them individually, i.e., we can work with real random variables normalized to variance one. This is the content of the following Lemma 6.5. Notice that the uniformity in the conditions (2.19) and (2.14) guarantees that the bounds (6.19) hold with uniform constants  $C_1, C_2$ . This implies the uniformity of the LSI constants, needed in Theorem 6.3, for the random variables constructed in Lemma 6.5. The proof of this Lemma will be given in Appendix C. We have thus proved Theorem 2.2.  $\square$

**Lemma 6.5** *Let  $m_3$  and  $m_4$  be two real numbers such that*

$$m_4 - m_3^2 - 1 \geq C_1, \quad m_4 \leq C_2 \quad (6.19)$$

*for some positive constants  $C_1$  and  $C_2$ . Then for any sufficient small  $\gamma > 0$  (depending on  $C_1$  and  $C_2$ ), there exists a real random variable  $\xi_\gamma$  whose distribution satisfies LSI and the first 4 moments of*

$$\xi' = (1 - \gamma)^{1/2} \xi_\gamma + \gamma^{1/2} \xi^G \quad (6.20)$$

*are 0, 1,  $m_3(\xi') = m_3$  and  $m_4(\xi')$ , and*

$$|m_4(\xi') - m_4| \leq C\gamma \quad (6.21)$$

*for some  $C$  depending on  $C_1$  and  $C_2$ , where  $\xi^G$  is real Gaussian random variable with mean 0 and variance 1, independent of  $\xi_\gamma$ . The LSI constant of  $\xi_\gamma$  (and thus  $\xi'$ ) is bounded from above by a function of  $C_1$  and  $C_2$ .*

## 7 Proof of Theorem 6.3

Theorem 6.3 states that the eigenvalues are at a distance  $N^{-1/2-\varepsilon}$  from their classical locations in a quadratic average sense. We will deduce this conclusion from the information on the closeness of the local density to the semicircle law. We note the constants appearing in this section may also depend on  $\alpha$  and  $\beta$  in (2.14), but we will not mention the dependence in the proof.

First we reformulate a result, which we have proved in [19], in a somewhat more general setup. It states that random points,  $\lambda_j$ , are close to a fixed set of locations,  $\gamma_j$ , if the local fluctuation is controlled, if the averaged counting function is close to the counting function of the  $\gamma_j$ 's in  $L^1$ -sense and if some tightness holds. For simplicity, the result is stated for the case when  $\gamma_j$ 's are the classical locations given by the semicircle law  $\varrho = \varrho_{sc}$ , (6.7), but the statement (and its proof) holds for any density function with support being a compact interval and with square root singularity at the edges. In particular, we applied this result in [19] for the Marchenko-Pastur (MP) distribution instead of the semicircle law. The counting function of  $\gamma_j$  can be replaced by its continuous version, i.e., by the distribution function of the semicircle law which defined by

$$n_{sc}(E) := \int_{-\infty}^E \varrho_{sc}(x) dx. \quad (7.1)$$

**Lemma 7.1** *Let  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N$  be an ordered collection of random points in  $\mathbb{R}$ . Denote the averaged counting function of  $\lambda_j$ 's*

$$n^\lambda(E) = \frac{1}{N} \mathbb{E} \#[\lambda_j \leq E]. \quad (7.2)$$

*Suppose the following four assumptions hold.*

1. *[Tightness at the edge] There exist  $m < 7$  and  $\varepsilon > 0$  such that*

$$n^\lambda(-2 - N^{-1/m}) \leq C e^{-N^\varepsilon} \quad \text{and} \quad n^\lambda(2 + N^{-1/m}) \geq 1 - C e^{-N^\varepsilon} \quad (7.3)$$

*and for any  $K \geq 3$ ,*

$$n^\lambda(K) \geq 1 - e^{-N^\varepsilon \log K} \quad \text{and} \quad n^\lambda(-K) \leq e^{-N^\varepsilon \log K}. \quad (7.4)$$

2. *[ $L^1$ -closeness of the counting functions]*

$$\int_{-\infty}^{\infty} |n^\lambda(E) - n_{sc}(E)| \, dE \leq C N^{-6/7}. \quad (7.5)$$

3. *[Fluctuation of moving averages] For any small  $\delta > 0$  there is a constant  $C$  such that for any  $j, K \in \mathbb{N}$  with  $j + K \leq N + 1$ , the local averages  $\lambda_{j,K} := K^{-1} \sum_{i=0}^{K-1} \lambda_{j+i}$  satisfy*

$$\mathbb{P} \left( |\lambda_{j,K} - \mathbb{E} \lambda_{j,K}| \geq N^{-1/2+\delta} K^{-1/2} \right) \leq C e^{-N^{\delta/2}}. \quad (7.6)$$

4. *[Positivity of the bulk density] There exists a small enough  $\delta > 0$  such that: for any interval  $I$  with  $|I| = N^{-5/8}$  and  $I \subset [-2 + N^{-\delta}, 2 - N^{-\delta}]$ , the number of the  $\lambda$ 's in  $I$  is bounded from below as follows*

$$\mathbb{P} \left( \#\{\lambda_j \in I\} \geq N^{-\delta} N |I| \right) \geq 1 - C N^{c \log \log N}. \quad (7.7)$$

*Then there exists  $\varepsilon > 0$  (independent of the constants in these four assumptions) such that*

$$\frac{1}{N} \sum_{i=1}^N \mathbb{E} (\lambda_i - \gamma_i)^2 \leq C N^{-1-\varepsilon} \quad (7.8)$$

*when  $N$  is large enough (depending on the constants in these four assumptions).*

*Proof of Lemma 7.1.* In Theorem 9.1 of [19] we have proved the analogous result on the singular values of the covariance matrix, where the role of the semicircle law was played by the MP law and the spectral edges,  $\pm 2$ , were replaced by  $\lambda_\pm$ , the two edges of the support of the MP distribution. In that paper we first proved the analogues of these four assumptions, then we presented the proof of (7.8) via a general argument that used only these assumptions. Inspecting the proofs of Lemma 9.5, 9.6 and 9.7 in [19], leading to (7.8), we observe that only equations (9.6), (9.8), (9.9) and (9.13) from [19] were used, in addition to the lower bound on the density of the points in the scale  $N^{-5/8}$ , which is used below (9.51) of [19]. The lower bound on the density is granted by the last assumption (7.7) (even with a better control on the probability than we required in [19]). Repeating the argument from [19], for the proof of Lemma 7.1 it is sufficient to check that the first three assumptions in Lemma 7.1 imply equations (9.6), (9.8), (9.9) and (9.13) in [19]. We now explain how to obtain these necessary bounds from our assumptions.



The first condition (7.3) corresponds to the input for Lemma 9.2 in [19], in particular, the analogue of (9.6) of [19],

$$-2 - N^{-1/m} \leq \mathbb{E}\lambda_j \leq 2 + N^{1/m},$$

follows immediately from (7.3) and (7.4). We note that (9.6) in [19] contains a threshold  $N^{-1/5}$  but actually in the proof we only needed it to be much less than  $N^{-1/7}$  (see (9.36)–(9.37) of [19] for the application of (9.6)).

The second condition (7.5) corresponds to Eq. (9.8) in [19]. As we showed in the proof of Lemma 9.3 of [19], Eq. (9.9) directly follows from (9.8). Here the analogous bound

$$\sup_E |n^\lambda(E) - n_{sc}(E)| \leq CN^{3/7}$$

follows directly from (7.5) in the same way.

Finally, the third condition (7.6) is exactly the same as (9.13) in [19]. Simply repeating now the proof of Theorem 9.1 from [19], we proved Lemma 7.1.  $\square$

Theorem 6.3 will now follow from Lemma 7.1 if we prove that the four conditions in the Lemma 7.1 hold in the case of generalized Wigner matrices (2.6). The last condition (7.7) follows from the local semicircle law (Theorem 2.1) and from the fact that  $\varrho_{sc}(x) \geq c\sqrt{\kappa}$  for  $x \in (-2 + \kappa, 2 - \kappa)$ . Here we list the first three conditions as three separate lemmas that will be proven in the next three subsections. This will complete the proof of Theorem 6.3.  $\square$

**Lemma 7.2** (1) *Let  $H$  be a generalized Wigner matrix with subexponential decay, in fact it is sufficient to assume that (2.14) and the upper bound  $C_{sup} < \infty$  in (2.6) hold. Define  $n^\lambda(E)$  as in (7.2). Then*

$$n^\lambda(-2 - N^{-1/6+\varepsilon}) \leq Ce^{-N^{\varepsilon'}} \quad \text{and} \quad n^\lambda(2 + N^{-1/6+\varepsilon}) \geq 1 - Ce^{-N^{\varepsilon'}} \quad (7.9)$$

for any small  $\varepsilon > 0$  with an  $\varepsilon' > 0$  depending on  $\varepsilon$ . Furthermore, for  $K \geq 3$ ,

$$n^\lambda(-K) \leq e^{-N^\varepsilon \log K} \quad \text{and} \quad n^\lambda(K) \geq 1 - e^{-N^\varepsilon \log K} \quad (7.10)$$

for some  $\varepsilon > 0$ .

(2) *In fact, the last tightness bound holds in a more general situation, namely, let the universal Wigner matrix  $H$  satisfy (2.1), (2.14) and  $M \geq (\log N)^9$  where  $M$  is defined in (2.3). Then we have*

$$n^\lambda(-3) \leq CN^{-c \log \log N} \quad \text{and} \quad n^\lambda(3) \geq 1 - CN^{-c \log \log N}. \quad (7.11)$$

**Lemma 7.3** *Let  $H$  satisfy the conditions of Theorem 6.3. Then for any  $\varepsilon > 0$  we have*

$$\int_{-\infty}^{\infty} |n^\lambda(E) - n_{sc}^\lambda(E)| dE \leq CN^{-1+\varepsilon}. \quad (7.12)$$

**Lemma 7.4** *Let  $H$  satisfy the conditions in Theorem 6.3, in particular, let the distribution of the matrix elements satisfy the uniform LSI (6.14). For  $j, K \in \mathbb{N}$ ,  $j + K \leq N + 1$ , define  $\lambda_{j,K} = K^{-1} \sum_{i=0}^{K-1} \lambda_{j+i}$ . Then for any  $\delta > 0$  small enough,*

$$\mathbb{P} \left( |\lambda_{j,K} - \mathbb{E}(\lambda_{j,K})| \geq N^{-1/2+\delta} K^{-1/2} \right) \leq Ce^{-N^\delta}, \quad (7.13)$$

with  $C$  depending on  $C_{sup}$  in (2.6) and  $C_S$  in (6.14).

## 7.1 Proof of Lemma 7.2.

Extreme eigenvalues are typically controlled by the moment method, evaluating  $\mathbb{E} \operatorname{Tr} H^k$  for large  $k$  using some graphical representation. Our proof follows the standard path, but since we were unable to find a reference that would apply precisely to our case, we include the proof for completeness. The main technical estimate (7.18) is borrowed from [34]. We remark that if we use the strongest result in [34], one can improve the exponent  $1/6$  to  $1/4$  in (7.9).

We start with the proof of (7.9) and (7.10) in the case of generalized Wigner matrices (see (2.6)). First we truncate the random variables. With the assumption of subexponential decay of  $h_{ij}$ , for any small  $\delta > 0$ , one can find a  $\widehat{h}_{ij}$  such that

$$\mathbb{P}(\widehat{h}_{ij} = h_{ij}) \geq 1 - e^{-N^{\varepsilon'}} \quad (7.14)$$

and

$$|\widehat{h}_{ij}| \leq N^{-1/2+\delta}, \quad \mathbb{E}(\widehat{h}_{ij}) = 0, \quad \mathbb{E}(|\widehat{h}_{ij}|^2) \leq \mathbb{E}(|h_{ij}|^2) \quad (7.15)$$

for some small number  $\varepsilon'$ , depending on  $\delta$ . Then we only need to bound the spectral norm of the new matrix  $\widehat{H} = (\widehat{h}_{ij})$ . To prove (7.9), it only remains to prove that, for some small  $\varepsilon' > 0$ ,

$$\mathbb{P}(\|\widehat{H}\| \geq 2 + N^{-1/6+\varepsilon}) \leq e^{-N^{\varepsilon'}}. \quad (7.16)$$

With

$$\mathbb{P}(\|\widehat{H}\| \geq 2 + N^{-1/6+\varepsilon}) \leq \frac{\mathbb{E} \operatorname{Tr} \widehat{H}^k}{(2 + N^{-1/6+\varepsilon})^k},$$

for any even  $k$ , (7.16) follows from

$$\mathbb{E} \operatorname{Tr} \widehat{H}^{k_0} \leq 2^{k_0+O(\log N)}, \quad (7.17)$$

with the choice of  $k_0 = N^{1/6-\delta/3}$  and  $\delta = 3\varepsilon/2$ , since  $\|\widehat{H}\|^k \leq \operatorname{Tr} \widehat{H}^k$  for even powers. The proof of (7.10) is analogous.

To estimate  $\mathbb{E}(\operatorname{Tr} \widehat{H}^k)$  for  $k \in \mathbb{N}$ , we start with introducing some notations and concepts on graphs.

Let  $p$  and  $k$  be given integers. We define the concept of *ordered closed walk* of  $k$  edges on an abstract ordered set  $A_p := \{a_1, a_2, \dots, a_p\}$  of  $p$  elements with the natural ordering  $a_1 < a_2 < \dots < a_p$ . An ordered closed walk on  $p$  vertices with  $k$  edge is determined by a sequence  $\underline{w} = (w_1, w_2, \dots, w_k)$  of the elements of  $A_p$  with the following properties:

- i) Along the walk, the fresh vertices from  $A_p$  are adjoined in increasing order, i.e.,  $\max_{j \leq m} w_j \leq \max_{j \leq m-1} w_j + 1$ .
- ii)  $\{w_1, w_2, \dots, w_k\} = A_p$ , i.e., all points of  $A_p$  are visited.
- iii) Let  $\Gamma(\underline{w})$  denote the undirected graph associated with  $\underline{w}$ , i.e., the vertex set of  $\Gamma(\underline{w})$  is  $A_p$ , the edges are given by  $(w_1, w_2), (w_2, w_3), \dots, (w_k, w_1)$ ; with multiple edges as well as self-loops ( $w_i = w_{i+1}$  for some  $i$ ) allowed. Then every edge of  $\Gamma$  appears at least twice.

Let  $\mathcal{W}(k, p)$  denote the set of ordered closed walks on  $p$  vertices with  $k$  edges. Their number was estimated in Lemma 2.1 of [34]

$$W(k, p) := |\mathcal{W}(k, p)| \leq \binom{k}{2p-2} p^{2(k-2p+2)} 2^{2p-2}. \quad (7.18)$$

This bound will be sufficient for the proof of (7.9) with exponent  $1/6 + \varepsilon$ . We remark that Lemma 4.1 of [34] gives a different bound on (7.18) that is better by essentially a factor  $[(k - 2p)/p]^{k-2p}$ . Applying this bound, one could improve the exponent in (7.9) to  $1/4 + \varepsilon$  but we will not pursue this improvement here.

We also need the concept of *labelling* the elements of  $A_p$  by the set  $\{1, 2, \dots, N\}$ . A labelling is given by a function  $\ell : A_p \rightarrow \{1, 2, \dots, N\}$  and we require that  $\ell$  be injective. The set of such labelling functions is denoted by  $\mathcal{L}(p, N)$ .

With these notations, we have the formula

$$\begin{aligned} \mathbb{E} \text{Tr} \widehat{H}^k &= \sum_{i_1, i_2, \dots, i_k=1}^N \mathbb{E} \widehat{h}_{i_1 i_2} \widehat{h}_{i_2 i_3} \dots \widehat{h}_{i_k i_1} \\ &= \sum_{p=1}^{k/2+1} \sum_{\underline{w} \in \mathcal{W}(k, p)} \sum_{\ell \in \mathcal{L}(p, N)} \mathbb{E} \widehat{h}_{\ell(w_1)\ell(w_2)} \widehat{h}_{\ell(w_2)\ell(w_3)} \dots \widehat{h}_{\ell(w_k)\ell(w_1)}. \end{aligned} \quad (7.19)$$

To verify this formula, for any given sequence  $i_1, i_2, \dots, i_k$  on the l.h.s., let  $p$  denote the number of different elements in this sequence and let the set  $A_p$  be identified with these different elements in the order of their appearance (i.e. for any  $m$  we let  $a_m := i_s$  for some  $s$  if  $i_s \neq i_t$ ,  $t < s$ , and  $i_s$  is the  $m$ -th freshest element among  $i_1, i_2, \dots, i_s$ , i.e.,  $|\{i_1, i_2, \dots, i_{s-1}\}| = m - 1$ ). Let  $w_1, w_2, \dots, w_k$  encode the sequence  $i_1, i_2, \dots, i_k$  with the new labels  $a_1, a_2, \dots, a_p$ . One may think of the walk,  $w_1, w_2, \dots, w_k$ , as the topological structure of the sequence  $(i_1, i_2, \dots, i_k)$  where the original labels from the set  $\{1, 2, \dots, N\}$  have been replaced by abstract labels, defined intrinsically from the repetition structure of  $(i_1, i_2, \dots, i_k)$ . Formula (7.19) is a resummation of all sequences  $(i_1, i_2, \dots, i_k)$  in terms of topological walks (first and second sum) and then reintroducing the original labelling with  $\{1, 2, \dots, N\}$  (third sum). Since the first moment of  $\widehat{h}_{ij}$  vanishes and different matrix elements are independent, all terms on the right hand side have zero expectation in which at least one factor  $\widehat{h}_{ij}$  appears only once. This justifies the requirement iii) in the definition of the ordered closed walks. The restriction  $p \leq k/2 + 1$  in the summation then comes from iii). This proves (7.19).

To compute the expectation on the r.h.s. of the (7.19), we need to introduce the concept of the *skeleton of the walk*. Given  $\underline{w} \in \mathcal{W}(k, p)$ , its skeleton  $S(\underline{w})$  is the undirected graph on  $A_p$  that is obtained from  $\Gamma(\underline{w})$  after replacing each multiple (parallel) edge by a single undirected edge. Here  $S(\underline{w})$  allows self-loops (as long as every edge has multiplicity 1). Thus the edge set  $E(S(\underline{w}))$  of the skeleton coincides with the edge set  $E(\Gamma(\underline{w}))$  after neglecting multiplicity and direction. The skeleton is a subgraph of the complete graph on  $A_p$ . We will also define the *tree of the walk*,  $T(\underline{w})$ , which is just a spanning tree of the skeleton  $S(\underline{w})$  built up successively along the walk by a greedy algorithm: include an edge to the  $T(\underline{w})$  if it does not create a loop together with the previously adjoined edges. Since  $\Gamma(\underline{w})$  is connected, and then so is  $S(\underline{w})$ , thus  $T(\underline{w})$  is indeed a tree on  $p$  vertices, in particular the number of its edges is

$$|E(T(\underline{w}))| = p - 1. \quad (7.20)$$

and  $S(\underline{w}) \setminus T(\underline{w})$  has total edge multiplicity less than  $k - 2(p - 1)$ .

For any edge  $e \in E(S(\underline{w}))$  of the skeleton, let  $\nu(e)$  denote the multiplicity of  $e$  in  $\Gamma(\underline{w})$  (edges with both orientations are taken into account). Clearly

$$\sum_{e \in E(S)} \nu(e) = k \quad (7.21)$$

for any skeleton graph  $S = S(\underline{w})$  for  $w \in \mathcal{W}(k, p)$ . Finally, for a given edge  $e = (a_\alpha, a_\beta)$  in a subgraph of  $A_p$  and for any labelling  $\ell \in \mathcal{L}(p, N)$ , we define the induced labelling of the edge  $e$  by  $\ell(e) = (\ell(a_\alpha), \ell(a_\beta))$ .

With these notations we have

$$\left| \mathbb{E} \widehat{h}_{\ell(w_1)\ell(w_2)} \widehat{h}_{\ell(w_2)\ell(w_3)} \cdots \widehat{h}_{\ell(w_k)\ell(w_1)} \right| \leq \prod_{e \in E(S(\underline{w}))} \mathbb{E} |\widehat{h}_{\ell(e)}|^{\nu(e)}.$$

Note that  $|\widehat{h}_{ij}| = |\widehat{h}_{ji}|$ , therefore there is no ambiguity in the notation  $|\widehat{h}_{\ell(e)}|$ . Since  $|\widehat{h}| \leq N^{-1/2+\delta}$  and  $\nu(e) \geq 2$ , we have

$$\mathbb{E} |\widehat{h}_{\ell(e)}|^{\nu(e)} \leq N^{(-1/2+\delta)(\nu(e)-2)} \sigma_{\ell(e)}^2, \quad (7.22)$$

or, alternatively,

$$\mathbb{E} |\widehat{h}_{\ell(e)}|^{\nu(e)} \leq N^{(-1/2+\delta)\nu(e)}. \quad (7.23)$$

We will use (7.22) for the edges of the tree,  $e \in E(T(\underline{w}))$ , and we use (7.23) for the remaining edges  $e \in E(S(\underline{w})) \setminus E(T(\underline{w}))$ . We can now estimate (7.19) using (7.21) and (7.20):

$$\begin{aligned} |\mathbb{E} \text{Tr} \widehat{H}^k| &\leq \sum_{p=1}^{k/2+1} \sum_{\underline{w} \in \mathcal{W}(k,p)} \sum_{\ell \in \mathcal{L}(p,N)} \prod_{e \in E(S(\underline{w}))} \mathbb{E} |\widehat{h}_{\ell(e)}|^{\nu(e)} \\ &\leq \sum_{p=1}^{k/2+1} \sum_{\underline{w} \in \mathcal{W}(k,p)} N^{(-1/2+\delta)(k-2(p-1))} \sum_{\ell \in \mathcal{L}(p,N)} \prod_{e \in E(T(\underline{w}))} \sigma_{\ell(e)}^2 \\ &\leq \sum_{p=1}^{k/2+1} \sum_{\underline{w} \in \mathcal{W}(k,p)} N^{1+(-1/2+\delta)(k-2(p-1))}. \end{aligned} \quad (7.24)$$

In the last step we used that

$$\sum_{\ell \in \mathcal{L}(p,N)} \prod_{e \in E(T)} \sigma_{\ell(e)}^2 \leq N.$$

holds for any tree  $T$ . This identity follows from successively summing up the labels for vertices with degree one in  $T$  by using the identity  $\sum_i \sigma_{i_j}^2 = 1$ .

Using (7.18), we obtain the bound

$$|\mathbb{E} \text{Tr} \widehat{H}^k| \leq \sum_{p=1}^{k/2+1} S(k,p), \quad (7.25)$$

with

$$S(k,p) := \binom{k}{2p-2} p^{2(k-2p+2)} 2^{2p-2} N^{1+(-1/2+\delta)(k-2(p-1))}. \quad (7.26)$$

It is easy to show that

$$S(k,p-1) \leq \frac{N^{2\delta} k^6}{4N} S(k,p). \quad (7.27)$$

Choosing  $k = N^{1/6-\delta/3}$ , we have  $S(k,p-1) \leq S(k,p)$ . Inserting this into (7.25), we obtain (7.17) and complete the proof.

Now we prove (7.11) with the same method. Similarly, with the assumption on the distribution of  $h_{ij}$ , one can find a  $\widehat{h}_{ij}$  such that

$$\mathbb{P}(\widehat{h}_{ij} = h_{ij}) \geq 1 - CN^{-c \log \log N} \quad (7.28)$$

and

$$|\widehat{h}_{ij}| \leq M^{-1/2}n, \quad \mathbb{E}(\widehat{h}_{ij}) = 0, \quad \mathbb{E}(|\widehat{h}_{ij}|^2) \leq \mathbb{E}(|h_{ij}|^2) + N^{-c \log \log N} \quad (7.29)$$

for  $n = (\log N)(\log \log N)$ . Here  $\widehat{h}_{ij}$  can be obtained by considering the cutoff random variables  $h_{ij} \mathbf{1}(|h_{ij}| \leq M^{-1/2}(\log N)(\log \log N))$  and then slightly modifying them to recover their zero expectation value.

We can again bound  $|\mathbb{E} \operatorname{Tr} \widehat{H}^k|$  as in (7.25) but with a slightly different  $S(k, p)$ ; instead of the factor  $N^{1+(-1/2+\delta)(k-2(p-1))}$  we will have  $N \cdot M^{(-1/2+\delta)(k-2(p-1))}$  in the definition (7.26). These modified  $S(k, p)$  numbers satisfy

$$S(k, p-1) \leq \frac{n^2 k^6}{4M} S(k, p) \quad (7.30)$$

and

$$S(k, k/2 + 1) = 2^k \cdot N. \quad (7.31)$$

Choosing  $k = n$ , we have  $\frac{n^2 k^6}{4M} < 1$ . Thus we obtain

$$|\mathbb{E} \operatorname{Tr} \widehat{H}^k| \leq 2^k \cdot 2nN, \quad (7.32)$$

which implies (7.11).  $\square$

## 7.2 Proof of Lemma 7.3.

First we show that the estimate on the expectation of  $m - m_{sc}$  is better than the estimate (2.16) on  $m - m_{sc}$  itself.

**Lemma 7.5** *Assume that the  $N \times N$  generalized Wigner matrix  $H$  (see (2.6)) satisfies (2.1), (2.4), (2.5) and (2.14),  $\mathbb{E} h_{ij} = 0$ , for any  $1 \leq i, j \leq N$  (i.e. the assumptions of Theorem 3.1 apart from (3.3) hold). Then we have, with some  $C > 0$ ,*

$$|\mathbb{E} m(z) - m_{sc}(z)| \leq \frac{(\log N)^C}{(N\eta)(\kappa + \eta)} \quad (7.33)$$

for any  $z = E + i\eta$ ,  $\eta > 0$ .

As a preparation to the proof, we need the following technical lemma that we state under more general conditions so that it is applicable for universal Wigner matrices.

**Lemma 7.6** *With the assumption of Theorem 3.1, suppose (3.3) holds, we have the estimate*

$$\left| \mathbb{E} m(z) + \frac{1}{\mathbb{E} m(z) + z} \right| \leq \frac{(\log N)^{C_0} \sqrt{\kappa + \eta}}{(M\eta)g^2(z)} \quad (7.34)$$

for some sufficiently large positive constants  $C_0$  (depending on  $\alpha, \beta$  in (2.14)).

*Proof of Lemma 7.6.* Recall the definitions of  $\Omega_z^o$ ,  $\Omega_z^d$  and  $\widehat{\Omega}_z^r$  in (3.17), (3.12), (3.13) and (3.14) and we define

$$\Omega_z \equiv \Omega_z^o \cap \Omega_z^d \cap \widehat{\Omega}_z^r. \quad (7.35)$$

With (3.34), we have

$$\mathbb{P}(\Omega_z) \geq 1 - CN^{-c \log \log N}. \quad (7.36)$$

The r.h.s. of (7.34) is larger than  $N^{-2}$ . Then with (7.36) and  $|m(z)| \leq \eta^{-1} \leq M$  (see (3.31)), we only need to prove

$$\left| \mathbb{E} \mathbf{1}(\Omega_z) m(z) - \frac{1}{\mathbb{E} \mathbf{1}(\Omega_z) m(z) + z} \right| \leq \frac{C(\log N)^C \sqrt{\kappa + \eta}}{(M\eta)g(z)^2}. \quad (7.37)$$

Taking the expectation of the self consistent equation (4.11) with (4.12), we obtain that

$$\mathbb{E}[\mathbf{1}(\Omega_z) \cdot G_{ii}] + \mathbb{E} \left[ \mathbf{1}(\Omega_z) \left( z + \sum_j \sigma_{ij}^2 G_{jj} + \Upsilon_i \right)^{-1} \right] = 0. \quad (7.38)$$

For simplicity, we define

$$A_i := \mathbb{E}[\mathbf{1}(\Omega_z) \cdot G_{ii}], \quad A := \sum_i A_i/N.$$

Together with (7.35) and (7.36), we have

$$|A - m_{sc}|, |A_j - A| \ll 1.$$

Then, similarly to (5.12), on the event  $\Omega_z$  we have

$$|z + A| - \left| \sum_j \sigma_{ij}^2 A_j - A \right| - |\Upsilon| > C,$$

by using  $|z + m_{sc}| \geq 1$  and that on the set  $\Omega_z$ ,  $\Upsilon$  is small. Therefore, we can expand (7.38) as

$$\begin{aligned} 0 = & A_i + \frac{1}{z + A} - \frac{\sum_j \sigma_{ij}^2 A_j - A}{(z + A)^2} - \frac{\mathbb{E} \mathbf{1}(\Omega_z) \Upsilon_i}{(z + A)^2} \\ & + O \left( \frac{\mathbb{E} \left[ \mathbf{1}(\Omega_z) \left| \sum_j \sigma_{ij}^2 G_{jj} - A \right|^2 \right]}{(z + A)^3} \right) + O \left( \frac{\mathbb{E} \left[ \mathbf{1}(\Omega_z) |\Upsilon_i|^2 \right]}{(z + A)^3} \right). \end{aligned} \quad (7.39)$$

Then summing up  $1 \leq i \leq N$ , we obtain that

$$\left| A + \frac{1}{z + A} \right| \leq C \max_i \left| \mathbb{E}[\mathbf{1}(\Omega_z) \Upsilon_i] \right| + C \max_i \mathbb{E} \left[ \mathbf{1}(\Omega_z) \left| \sum_j \sigma_{ij}^2 G_{jj} - A \right|^2 \right] + C \mathbb{E}[\mathbf{1}(\Omega_z) |\Upsilon|^2]. \quad (7.40)$$

Applying (3.12) and the definition of  $\Omega_z$ , we can bound the second and third terms in the r.h.s. of (7.40) with some constant  $C$  as follows,

$$\left| A + \frac{1}{z + A} \right| \leq C \max_i \left| \mathbb{E}[\mathbf{1}(\Omega_z) \Upsilon_i] \right| + \frac{(\log N)^C \sqrt{\kappa + \eta}}{M\eta g(z)^2}. \quad (7.41)$$

If  $\eta > 3$ , we estimate  $\mathbb{E}[\mathbf{1}(\Omega_z) \Upsilon_i]$  as

$$\left| \mathbb{E}[\mathbf{1}(\Omega_z) \Upsilon_i] \right| \leq \left| \mathbb{E}[\mathbf{1}(\widehat{\Omega}_z^o) \Upsilon_i] \right| + \mathbb{E} \left[ \mathbf{1}([\widehat{\Omega}_z^o]^c) \mathbf{1}(\widehat{\Omega}_z) |\Upsilon_i| \right]. \quad (7.42)$$

With (4.7), we have

$$|G_{ij} G_{ji} / G_{ii}| \leq 2/\eta. \quad (7.43)$$

Then, with the definition of  $\Upsilon_i$  (4.12) and (4.39), we have

$$\mathbb{P}(|\max \Upsilon_i| \geq N^C) \leq e^{-N^c} \quad (7.44)$$

for some positive constants  $c$  and  $C$ . Inserting this and (3.20) into (7.42), we have

$$\left| A + \frac{1}{z + A} \right| \leq \frac{(\log N)^C \sqrt{\kappa + \eta}}{M\eta g^2(z)} \quad (7.45)$$

in the case of  $\eta > 3$ . If  $\eta < 3$ , similarly, with (3.23) we have the same result. This proves (7.37) and thus completes the proof of Lemma 7.6.  $\square$

*Proof of Lemma 7.5.* First we will prove the result for large  $\eta$ , more precisely we show (7.33) under the additional assumption that

$$N\eta(\kappa + \eta)^{3/2} \geq (\log N)^{C_1}, \quad (7.46)$$

with a sufficiently large constant  $C_1$ .

In the case of the generalized Wigner matrix, (2.6), we have  $M \geq (C_{sup})^{-1}N$  and  $\delta_+ \geq C_{inf}$  (2.7), then

$$g(z) \asymp \sqrt{\kappa + \eta}$$

up to an  $O(1)$  factor. Note that with a sufficiently large  $C_1$ , (7.46) implies (3.3) and thus combining Lemma 7.6 with Lemma 5.2 we obtain (7.33) under the condition that  $\eta$  satisfies (7.46).

To prove (7.33) for any  $\eta > 0$ , it remains to consider the case when (7.46) does not hold. For a fixed  $E$ , let  $\eta^* = \eta^*(E) > 0$  be the (unique) solution of  $N\eta(\kappa + \eta)^{3/2} = (\log N)^{C_1}$ , i.e. when (7.46) becomes an equality. In particular, we know that

$$|\mathbb{E}m(z^*) - m_{sc}(z^*)| \leq \frac{(\log N)^C}{(N\eta^*)(\kappa + \eta^*)}. \quad (7.47)$$

Consider  $\eta < \eta^*$ , set  $z = E + i\eta$ ,  $z^* = E + i\eta^*$  and estimate

$$|\mathbb{E}m(z) - m_{sc}(z)| \leq |\mathbb{E}m(z^*) - m_{sc}(z^*)| + \int_{\eta}^{\eta^*} |\partial_y(\mathbb{E}m(E + iy) - m_{sc}(E + iy))| dy. \quad (7.48)$$

Note that

$$|\partial_y m(E + iy)| = \left| \frac{1}{N} \sum_j \partial_y G_{jj}(E + iy) \right| \quad (7.49)$$

$$\leq \frac{1}{N} \sum_{jk} |G_{jk}(E + iy)|^2 = \frac{1}{Ny} \sum_j \text{Im} G_{jj}(E + iy) = \frac{1}{y} \text{Imm}(E + iy), \quad (7.50)$$

and similarly

$$|\partial_y m_{sc}(E + iy)| = \left| \int \frac{\varrho_{sc}(x)}{(x - E - iy)^2} dx \right| \leq \int \frac{\varrho_{sc}(x)}{|x - E - iy|^2} dx = \frac{1}{y} \text{Imm}_{sc}(E + iy).$$

Now we use the fact that the functions  $y \rightarrow y\mathbb{I}m(E + iy)$  and  $y \rightarrow y\mathbb{I}m_{sc}(E + iy)$  are monotone increasing for any  $y > 0$  since both are Stieltjes transforms of a positive measure. Therefore the integral in (7.48) can be bounded by

$$\int_{\eta}^{\eta^*} \frac{dy}{y} [\mathbb{I}m\mathbb{E}m(E + iy) + \mathbb{I}m_{sc}(E + iy)] \leq \eta^* [\mathbb{I}m\mathbb{E}m(E + i\eta^*) + \mathbb{I}m_{sc}(E + i\eta^*)] \int_{\eta}^{\eta^*} \frac{dy}{y^2} \quad (7.51)$$

By the choice of  $\eta^*$  and using that  $\mathbb{I}m_{sc}(z^*) \leq C\sqrt{\kappa + \eta^*}$ , we have

$$\mathbb{I}m_{sc}(z^*) \leq \frac{(\log N)^C}{(N\eta^*)(\kappa + \eta^*)}. \quad (7.52)$$

and then  $\mathbb{I}m\mathbb{E}m(z^*)$  can be estimated from (7.47). Inserting these estimates into (7.48) and (7.51), and using (7.47), we get

$$|\mathbb{E}m(z) - m_{sc}(z)| \leq |\mathbb{E}m(z^*) - m_{sc}(z^*)| + \frac{2(\log N)^C}{N\eta^*(\kappa + \eta^*)} \frac{\eta^*}{\eta} \leq \frac{(\log N)^C}{N\eta(\kappa + \eta)}$$

with a possible larger  $C$  in the r.h.s. This completes the proof of Lemma 7.5.  $\square$

With Lemma 7.5, it follows that for any  $E$  and  $\eta > 0$ ,

$$|n^\lambda(E + \eta) - n^\lambda(E - \eta)| + |n_{sc}^\lambda(E + \eta) - n_{sc}^\lambda(E - \eta)| \leq \eta(\log N)^C \left(1 + \frac{1}{N\eta(|E - 2| + \eta)}\right). \quad (7.53)$$

Now we return to the main argument to prove (7.12) in Lemma 7.3. Given (7.10), we only need to prove

$$\int_{-3}^3 |n^\lambda(E) - n_{sc}^\lambda(E)| dE \leq CN^{-1+\varepsilon}. \quad (7.54)$$

This inequality follows from the next lemma by choosing the signed measure

$$\varrho^\Delta(dx) = \varrho_{sc}(dx) - \frac{dn^\lambda(E)}{dE}, \quad (7.55)$$

whose Stieltjes transform is given by

$$m^\Delta(z) = m_{sc}(z) - \mathbb{E}m(z) \quad (7.56)$$

and the conditions (7.58) and (7.59) are provided by (7.33) and (7.53). This will complete the proof of Lemma 7.3.

**Lemma 7.7** *Let  $\varrho^\Delta(dx)$  be a finite signed measure with support in  $[-K, K]$  for some  $K > 0$ . Let*

$$m^\Delta(z) := \int_{\mathbb{R}} \frac{\varrho^\Delta(dx)}{x - z}, \quad n^\Delta(E) := \int_{-\infty}^E \varrho^\Delta(dx) \quad (7.57)$$

*be the Stieltjes transform and the distribution function of  $\varrho^\Delta(dx)$ , respectively. Let  $\kappa_x, \kappa_E$  denote  $\||x| - 2|$  and  $\||E| - 2|$ . We assume that  $m^\Delta$  satisfies the following bound with some constant  $C$ :*

$$|m^\Delta(x + iy)| \leq \frac{(\log N)^C}{(Ny)(\kappa_x + y)} \quad \text{for } y > 0, \quad |x| \leq K + 1, \quad (7.58)$$



and for any  $a > 0$

$$\int_{E-a}^{E+a} |\varrho^\Delta|(dx) \leq a(\log N)^C \left(1 + \frac{1}{Na(\kappa_E + a)}\right). \quad (7.59)$$

Then

$$\int_{-K}^K dE |n^\Delta(E)| \leq CN^{-1}(\log N)^{C'} \quad (7.60)$$

for some constant  $C' > 0$  when  $N$  is sufficiently large.

This lemma is similar to Lemma B.1 in [16], but with different assumptions. Since the assumptions here are stronger than (B.3) and (B.4) in [16], we actually obtain a better bound (7.60) than in [16], where the l.h.s. of (7.60) was bounded by  $N^{-6/7}$ .

*Proof of Lemma 7.7.* For simplicity, we omit the  $\Delta$  superscript in the proof. For a fixed  $E \in [-K, K]$ ,  $\eta > 0$ , define a function  $f = f_{E,\eta}: \mathbb{R} \rightarrow \mathbb{R}$ : such that  $f(x) = 1$  for  $x \in [-K, E - \eta]$ ,  $f(x)$  vanishes for  $x \in (-\infty, -K - 1) \cap [E + \eta, \infty)$ , moreover  $|f'(x)| \leq C\eta^{-1}$  and  $|f''(x)| \leq C\eta^{-2}$ . Then

$$\left| n(E) - \int_{\mathbb{R}} f_{E,\eta}(\lambda) \varrho(\lambda) d\lambda \right| \leq \int_{E-\eta}^{E+\eta} |\varrho|(dx) \leq \eta(\log N)^C \left(1 + \frac{1}{N\eta(\kappa_E + \eta)}\right). \quad (7.61)$$

We will choose  $\eta = N^{-1}$  and set  $f_E := f_{E,\eta}$  with  $\eta = 1/N$ . Then to prove (7.60), we only need to prove that

$$\left| \int_{|E| \leq K+1} \int_{\mathbb{R}} f_E(\lambda) \varrho(\lambda) d\lambda dE \right| \leq N^{-1}(\log N)^{C'} \quad (7.62)$$

for some  $C' > 0$ .

To express  $f_E(\lambda)$  in terms of the Stieltjes transform, we use the Helffer-Sjöstrand functional calculus, as (B.12) in [16]. We formulate this result in a more general form.

**Lemma 7.8** *Let  $f_{E,\eta}$  be given as above with some  $E \in [-K, K]$ ,  $K \geq 3$ , and  $0 < \eta \leq 1/2$ . Suppose that the Stieltjes transform  $m$  of the signed measure  $\varrho$  satisfies*

$$|m(x + iy)| \leq \frac{L}{(Ny)^\tau (\kappa_x + y)^\sigma} \quad \text{for } y > 0, \quad |x| \leq K + 1, \quad (7.63)$$

with some exponents  $0 \leq \tau, \sigma \leq 1$  and some constant  $L$ . Then

$$\left| \int f_E(\lambda) \varrho(\lambda) d\lambda \right| \leq \frac{CL |\log \eta|}{N^\tau (\kappa_E + \eta)^\sigma}, \quad (7.64)$$

with some constant  $C$  depending on  $K$ .

The condition of this lemma with  $\tau = \sigma = 1$  and  $L = (\log N)^C$  coincides with (7.58), therefore, after integrating in  $E$  and using  $\eta = 1/N$ , we obtain (7.62) which completes the proof of Lemma 7.7.  $\square$

*Proof of Lemma 7.8.* Analogously to (B.13), (B.14) and (B.15) in [16] we obtain that

$$\begin{aligned}
\left| \int f_E(\lambda) \varrho(\lambda) d\lambda \right| &\leq C \int_{\mathbb{R}^2} (|f_E(x)| + |y| |f'_E(x)|) |\chi'(y)| |m(x+iy)| dx dy \\
&+ C \left| \int_{|y| \leq \eta} \int y f''_E(x) \chi(y) \operatorname{Im} m(x+iy) dx dy \right| \\
&+ C \left| \int_{|y| \geq \eta} \int_{\mathbb{R}} y f''_E(x) \chi(y) \operatorname{Im} m(x+iy) dx dy \right|,
\end{aligned} \tag{7.65}$$

where  $\chi(y)$  is a smooth cutoff function with support in  $[-1, 1]$ , with  $\chi(y) = 1$  for  $|y| \leq 1/2$  and with bounded derivatives. The first term is estimated by

$$\int_{\mathbb{R}^2} (|f_E(x)| + |y| |f'_E(x)|) |\chi'(y)| |m(x+iy)| dx dy \leq \frac{CL}{N^\tau}, \tag{7.66}$$

using (7.63) and the support of  $\chi'$ .

With (7.63) and  $|f''_E| \leq C\eta^{-2}$  and

$$\operatorname{supp} f'_E(x) \subset \{|x - E| \leq \eta\},$$

the second term in r.h.s. of (7.65) is bounded by

$$\begin{aligned}
CL \left| \int_{0 \leq y \leq \eta} \int_{|x-E| \leq \eta} \frac{y |f''_E(x)|}{(Ny)^\tau (\kappa_x + y)^\sigma} dx dy \right| &\leq \frac{CL}{N^\tau \eta^2} \left| \int_{0 \leq y \leq \eta} \int_{|x-E| \leq \eta} \frac{y^{1-\tau}}{(\kappa_x + y)^\sigma} dx dy \right| \\
&\leq \frac{CL \eta^{1-\tau} |\log \eta|}{N^\tau (\kappa_E + \eta)^\sigma}.
\end{aligned} \tag{7.67}$$

Here we used that for  $y \leq 1/2$  we have

$$\int_{|x-E| \leq \eta} \frac{1}{(\kappa_x + y)^\sigma} dx \leq \frac{C\eta |\log y|}{(\kappa_E + \eta)^\sigma}.$$

As the (B.17) and (B.19) in [16], we integrate the third term in (7.65) by parts first in  $x$ , then in  $y$ . Then bound it with absolute value by

$$C \int_{|x| \leq K+1} \eta |f'_E(x)| |\operatorname{Re} m(x+i\eta)| dx + C \int_{\mathbb{R}^2} |f'_E(x) \chi'(y) \operatorname{Re} m(x+iy)| + \frac{C}{\eta} \int_{\eta \leq y \leq 1} \int_{|x-E| \leq \eta} |\operatorname{Re} m(x+iy)| dx dy. \tag{7.68}$$

The middle term is bounded as (7.66). With (7.63) again, we have

$$\begin{aligned}
(7.68) &\leq \frac{CL}{(N\eta)^\tau} \int_{|x-E| \leq \eta} \frac{1}{(\kappa_x + \eta)^\sigma} dx + \frac{CL}{(N\eta)^\tau} + \frac{CL}{(N\eta)^\tau} \int_{\eta \leq y \leq 1} \int_{|x-E| \leq \eta} \frac{1}{(\kappa_x + y)^\sigma} dx dy \\
&\leq \frac{CL \eta^{1-\tau} |\log \eta|}{N^\tau (\kappa_E + \eta)^\sigma}.
\end{aligned} \tag{7.69}$$

Then combining (7.65), (7.66), (7.67), (7.68) and (7.69) we obtain (7.64) and complete the proof of Lemma 7.8.  $\square$

### 7.3 Proof of Lemma 7.4

Define the variables  $v_{ij}$  as

$$h_{ij} = \sigma_{ij} v_{ij}. \quad (7.70)$$

Denote by  $u_\alpha$  and  $\lambda_\alpha$  the eigenvectors and eigenvalues of  $H$ . For any collection of real numbers,  $C_\alpha \in \mathbb{R}$ , we have

$$\sum_{ij} \left| \sum_{\alpha} C_{\alpha} \frac{\partial \lambda_{\alpha}}{\partial v_{ij}} \right|^2 = \sum_{ij} \left| \sum_{\alpha} C_{\alpha} \sigma_{ij} \bar{u}_{\alpha}(i) u_{\alpha}(j) \right|^2 = \sum_{ij} \sigma_{ij}^2 \left| \sum_{\alpha} C_{\alpha} \bar{u}_{\alpha}(i) u_{\alpha}(j) \right|^2 \leq C_{sup} N^{-1} \sum_{\alpha} |C_{\alpha}|^2. \quad (7.71)$$

With the choice  $C_{\alpha} = K^{-1}$ ,  $\alpha = j, j+1, \dots, j+K-1$ , and  $C_{\alpha} = 0$  otherwise, we get  $|\nabla \lambda_{j,K}|^2 \leq C_{sup} (NK)^{-1}$ . Using the Bobkov-Götze concentration inequality [4] and the uniform bound on the LSI constant (6.14), we get

$$\mathbb{P}(|\lambda_{j,K} - \mathbb{E} \lambda_{j,K}| \geq \gamma) \leq e^{-\gamma T} \mathbb{E} e^{C_S T^2 |\nabla \lambda_{j,K}|^2} \leq e^{-\gamma T + C_S C_{sup} T^2 / (NK)}$$

for any  $T$  and  $\gamma$ . Choosing  $\gamma = N^{-1/2+\delta} K^{-1/2}$  and  $T = (NK)^{1/2}$ , we obtain (7.13).  $\square$

## 8 Proof of the Green's function comparison theorem

*Proof of Theorem 2.3.* From the trivial bound

$$\operatorname{Im} \left( \frac{1}{H - E - i\eta} \right)_{jj} \leq \left( \frac{y}{\eta} \right) \operatorname{Im} \left( \frac{1}{H - E - iy} \right)_{jj}, \quad \eta \leq y,$$

and from (2.21) we have the following a priori bound

$$\mathbb{P} \left( \max_{0 \leq \gamma \leq \gamma(N)} \max_{1 \leq k \leq N} \max_{|E| \leq 2^{-\kappa}} \sup_{\eta \geq N^{-1-\epsilon}} \left| \operatorname{Im} \left( \frac{1}{H_{\gamma} - E \pm i\eta} \right)_{kk} \right| \leq N^{3\tau+\epsilon} \right) \geq 1 - CN^{-c \log \log N}. \quad (8.1)$$

Note that the supremum over  $\eta$  can be included by establishing the estimate first for a fine grid of  $\eta$ 's with spacing  $N^{-10}$  and then extend the bound for all  $\eta$  by using that the Green's functions are Lipschitz continuous in  $\eta$  with a Lipschitz constant  $\eta^{-2}$ .

Let  $\lambda_m$  and  $u_m$  denote the eigenvalues and eigenvectors of  $H_{\gamma}$ , then by the definition of the Green's function, we have

$$\left| \left( \frac{1}{H_{\gamma} - z} \right)_{jk} \right| \leq \sum_{m=1}^N \frac{|u_m(j)| |u_m(k)|}{|\lambda_m - z|} \leq \left[ \sum_{m=1}^N \frac{|u_m(j)|^2}{|\lambda_m - z|} \right]^{1/2} \left[ \sum_{m=1}^N \frac{|u_m(k)|^2}{|\lambda_m - z|} \right]^{1/2}.$$

Define a dyadic decomposition

$$U_n = \{m : 2^{n-1}\eta \leq |\lambda_m - E| < 2^n \eta\}, \quad n = 1, 2, \dots, n_0 := C \log N, \quad (8.2)$$

$$U_0 = \{m : |\lambda_m - E| < \eta\}, \quad U_{\infty} := \{m : 2^{n_0} \eta \leq |\lambda_m - E|\},$$

and divide the summation over  $m$  into  $\cup_n U_n$

$$\sum_{m=1}^N \frac{|u_m(j)|^2}{|\lambda_m - z|} = \sum_n \sum_{m \in U_n} \frac{|u_m(j)|^2}{|\lambda_m - z|} \leq C \sum_n \sum_{m \in U_n} \operatorname{Im} \frac{|u_m(j)|^2}{\lambda_m - E - i2^n \eta} \leq C \sum_n \operatorname{Im} \left( \frac{1}{H_{\gamma} - E - i2^n \eta} \right)_{jj}.$$

Using the estimate (2.21) for  $n = 0, 1, \dots, n_0$  and a trivial bound of  $O(1)$  for  $n = \infty$ , we have proved that

$$\mathbb{P} \left( \sup_{0 \leq \gamma \leq \gamma(N)} \sup_{1 \leq k, \ell \leq N} \max_{|E| \leq 2^{-\kappa}} \sup_{\eta \geq N^{-1-\varepsilon}} \left| \left( \frac{1}{H_\gamma - E \pm i\eta} \right)_{k\ell} \right| \leq N^{4\tau+\varepsilon} \right) \geq 1 - CN^{-c \log \log N}. \quad (8.3)$$

For simplicity, we will consider the case when the test function  $F$  has only  $n = 1$  variable and  $k_1 = 1$ , i.e., we consider the trace of a first order monomial; the general case follows analogously. Consider the telescoping sum of differences of expectations

$$\begin{aligned} \mathbb{E} F \left( \frac{1}{N} \text{Tr} \frac{1}{H^{(v)} - z} \right) - \mathbb{E} F \left( \frac{1}{N} \text{Tr} \frac{1}{H^{(w)} - z} \right) \\ = \sum_{\gamma=1}^{\gamma(N)} \left[ \mathbb{E} F \left( \frac{1}{N} \text{Tr} \frac{1}{H_\gamma - z} \right) - \mathbb{E} F \left( \frac{1}{N} \text{Tr} \frac{1}{H_{\gamma-1} - z} \right) \right]. \end{aligned} \quad (8.4)$$

Let  $E^{(ij)}$  denote the matrix whose matrix elements are zero everywhere except at the  $(i, j)$  position, where it is 1, i.e.,  $E_{k\ell}^{(ij)} = \delta_{ik}\delta_{j\ell}$ . Fix an  $\gamma \geq 1$  and let  $(i, j)$  be determined by  $\phi(i, j) = \gamma$ . We will compare  $H_{\gamma-1}$  with  $H_\gamma$ . Note that these two matrices differ only in the  $(i, j)$  and  $(j, i)$  matrix elements and they can be written as

$$\begin{aligned} H_{\gamma-1} &= Q + \frac{1}{\sqrt{N}} V, & V &:= v_{ij} E^{(ij)} + v_{ji} E^{(ji)} \\ H_\gamma &= Q + \frac{1}{\sqrt{N}} W, & W &:= w_{ij} E^{(ij)} + w_{ji} E^{(ji)}, \end{aligned}$$

with a matrix  $Q$  that has zero matrix element at the  $(i, j)$  and  $(j, i)$  positions and where we set  $v_{ji} := \bar{v}_{ij}$  for  $i < j$  and similarly for  $w$ . Define the Green's functions

$$R = \frac{1}{Q - z}, \quad S = \frac{1}{H_\gamma - z}.$$

We first claim that the estimate (8.3) holds for the Green's function  $R$  as well. To see this, we have, from the resolvent expansion,

$$R = S + N^{-1/2} S V S + \dots + N^{-9/5} (S V)^9 S + N^{-5} (S V)^{10} R.$$

Since  $V$  has only at most two nonzero element, when computing the  $(k, \ell)$  matrix element of this matrix identity, each term is a finite sum involving matrix elements of  $S$  or  $R$  and  $v_{ij}$ , e.g.  $(S V S)_{k\ell} = S_{ki} v_{ij} S_{j\ell} + S_{kj} v_{ji} S_{i\ell}$ . Using the bound (8.3) for the  $S$  matrix elements, the subexponential decay for  $v_{ij}$  and the trivial bound  $|R_{ij}| \leq \eta^{-1}$ , we obtain that the estimate (8.3) holds for  $R$ .

We can now start proving the main result. By the resolvent expansion,

$$S = R - N^{-1/2} R V R + N^{-1} (R V)^2 R - N^{-3/2} (R V)^3 R + N^{-2} (R V)^4 R - N^{-5/2} (R V)^5 S,$$

so we can write

$$\frac{1}{N} \text{Tr} S = \widehat{R} + \xi, \quad \xi = \sum_{m=1}^4 N^{-m/2} \widehat{R}^{(m)} + N^{-5/2} \Omega$$

with

$$\widehat{R} = \frac{1}{N} \text{Tr } R, \quad \widehat{R}^{(m)} = (-1)^m \frac{1}{N} \text{Tr } (RV)^m R, \quad \Omega = -\frac{1}{N} \text{Tr } (RV)^5 S.$$

For each diagonal element in the computation of these traces, the contribution to  $\widehat{R}$ ,  $\widehat{R}^{(m)}$  and  $\Omega$  is a sum of a few terms. E.g.

$$\widehat{R}^{(2)} = \frac{1}{N} \sum_k \left[ R_{ki} v_{ij} R_{jj} v_{ji} R_{ik} + R_{ki} v_{ij} R_{ji} v_{ij} R_{jk} + R_{kj} v_{ji} R_{ii} v_{ij} R_{jk} + R_{kj} v_{ji} R_{ij} v_{ji} R_{ik} \right]$$

and similar formulas hold for the other terms.

Then we have

$$\begin{aligned} \mathbb{E} F \left( \frac{1}{N} \text{Tr} \frac{1}{H_\gamma - z} \right) &= \mathbb{E} F \left( \widehat{R} + \xi \right) \\ &= \mathbb{E} \left[ F(\widehat{R}) + F'(\widehat{R})\xi + F''(\widehat{R})\xi^2 + \dots + F^{(5)}(\widehat{R} + \xi')\xi^5 \right] \\ &= \sum_{m=0}^5 N^{-m/2} \mathbb{E} A^{(m)}, \end{aligned} \tag{8.5}$$

where  $\xi'$  is a number between 0 and  $\xi$  and it depends on  $\widehat{R}$  and  $\xi$ ; the  $A^{(m)}$ 's are defined as

$$A^{(0)} = F(\widehat{R}), \quad A^{(1)} = F'(\widehat{R})\widehat{R}^{(1)}, \quad A^{(2)} = F''(\widehat{R})(\widehat{R}^{(1)})^2 + F'(\widehat{R})\widehat{R}^{(2)},$$

and similarly for  $A^{(3)}$  and  $A^{(4)}$ . Finally,

$$A^{(5)} = F'(\widehat{R})\Omega + F^{(5)}(\widehat{R} + \xi')(\widehat{R}^{(1)})^5 + \dots$$

The expectation values of the terms  $A^{(m)}$ ,  $m \leq 4$ , with respect to  $v_{ij}$  are determined by the first four moments of  $v_{ij}$ , for example

$$\begin{aligned} \mathbb{E} A^{(2)} &= F'(\widehat{R}) \left[ \frac{1}{N} \sum_k R_{ki} R_{jj} R_{ik} + \dots \right] \mathbb{E} |v_{ij}|^2 + F''(\widehat{R}) \left[ \frac{1}{N^2} \sum_{k,\ell} R_{ki} R_{j\ell} R_{\ell j} R_{ik} + \dots \right] \mathbb{E} |v_{ij}|^2 \\ &\quad + F'(\widehat{R}) \left[ \frac{1}{N} \sum_k R_{ki} R_{ji} R_{jk} + \dots \right] \mathbb{E} v_{ij}^2 + F''(\widehat{R}) \left[ \frac{1}{N^2} \sum_{k,\ell} R_{ki} R_{j\ell} R_{\ell i} R_{jk} + \dots \right] \mathbb{E} v_{ij}^2. \end{aligned}$$

Note that the coefficients involve up to four derivatives of  $F$  and normalized sums of matrix elements of  $R$ . Using the estimate (8.3) for  $R$  and the derivative bounds (2.23) for the typical values of  $\widehat{R}$ , we see that all these coefficients are bounded by  $N^{C(\tau+\varepsilon)}$  with a very large probability, where  $C$  is an explicit constant. We use the bound (2.24) for the extreme values of  $\widehat{R}$  but this event has a very small probability by (8.3). Therefore, the coefficients of the moments  $\mathbb{E} \bar{v}_{ij}^s v_{ij}^u$ ,  $u + s \leq 4$ , in the quantities  $A^{(0)}, \dots, A^{(4)}$  are essentially bounded, modulo a factor  $N^{C(\tau+\varepsilon)}$ . Notice that the fourth moment of  $v_{ij}$  appears only in the  $m = 4$  term that already has a prefactor  $N^{-2}$  in (8.5). Therefore, to compute the  $m \leq 4$  terms in (8.5) up to a precision  $o(N^{-2})$ , it is sufficient to know the first three moments of  $v_{ij}$  exactly and the fourth moment only with a precision  $N^{-\delta}$ ; if  $\tau$  and  $\varepsilon$  are chosen such that  $C(\tau + \varepsilon) < \delta$ , then the discrepancy in the fourth moment is irrelevant.

Finally, we have to estimate the error term  $A^{(5)}$ . All terms without  $\Omega$  can be dealt with as before; after estimating the derivatives of  $F$  by  $N^{C(\tau+\varepsilon)}$ , one can perform the expectation with respect to  $v_{ij}$  that is independent of  $\widehat{R}^{(m)}$ . For the terms involving  $\Omega$  one can argue similarly, by appealing to the fact that the matrix elements of  $S$  are also essentially bounded by  $N^{C(\tau+\varepsilon)}$ , see (8.3), and that  $v_{ij}$  has subexponential decay. Alternatively, one can use Hölder inequality to decouple  $S$  from the rest and use (8.3) directly, for example:

$$\mathbb{E}|F'(\widehat{R})\Omega| = \frac{1}{N}\mathbb{E}|F'(\widehat{R})\text{Tr}(RV)^5S| \leq \frac{1}{N}\left[\mathbb{E}(F'(\widehat{R}))^2\text{Tr}S^2\right]^{1/2}\left[\mathbb{E}\text{Tr}(RV)^5(VR^*)^5\right]^{1/2} \leq CN^{C(\tau+\varepsilon)}.$$

Note that exactly the same perturbation expansion holds for the resolvent of  $H_{\gamma-1}$ , just  $v_{ij}$  is replaced with  $w_{ij}$  everywhere. By the moment matching condition, the expectation values  $\mathbb{E}A^{(m)}$  of terms for  $m \leq 3$  in (8.5) are identical and the  $m = 4$  term differs by  $N^{-\delta+C(\tau+\varepsilon)}$ . Choosing  $\tau = \varepsilon$ , we have

$$\mathbb{E}F\left(\frac{1}{N}\text{Tr}\frac{1}{H_\gamma - z}\right) - \mathbb{E}F\left(\frac{1}{N}\text{Tr}\frac{1}{H_{\gamma-1} - z}\right) \leq CN^{-5/2+C\varepsilon} + CN^{-2-\delta+C\varepsilon}.$$

After summing up in (8.4) we have thus proved that

$$\mathbb{E}F\left(\frac{1}{N}\text{Tr}\frac{1}{H^{(v)} - z}\right) - \mathbb{E}F\left(\frac{1}{N}\text{Tr}\frac{1}{H^{(w)} - z}\right) \leq CN^{-1/2+C\varepsilon} + CN^{-\delta+C\varepsilon}.$$

The proof can be easily generalized to functions of several variables. This concludes the proof of Theorem 2.3.  $\square$

*Proof of Theorem 6.4.* Define an approximate delta function (times  $\pi$ ) at the scale  $\eta$  by

$$\theta_\eta(x) = \text{Im}\frac{1}{x - i\eta}.$$

For notational simplicity, we will prove only the case of three point correlation functions; the proof is analogous for the general case. By definition of the correlation function, for any fixed  $E$ ,  $\alpha_1, \alpha_2, \alpha_3$ ,

$$\begin{aligned} & \mathbb{E}_{\mathbf{w}}\frac{1}{N(N-1)(N-2)}\sum_{i \neq j \neq k}\theta_\eta\left(\lambda_i - E - \frac{\alpha_1}{N}\right)\theta_\eta\left(\lambda_j - E - \frac{\alpha_2}{N}\right)\theta_\eta\left(\lambda_k - E - \frac{\alpha_3}{N}\right) \\ &= \int dx_1 dx_2 dx_3 p_{w,N}^{(3)}(x_1, x_2, x_3)\theta_\eta(x_1 - E_1)\theta_\eta(x_2 - E_2)\theta_\eta(x_3 - E_3), \quad E_j := E + \frac{\alpha_j}{N}. \end{aligned} \quad (8.6)$$

By the exclusion-inclusion principle,

$$\mathbb{E}_{\mathbf{w}}\frac{1}{N(N-1)(N-2)}\sum_{i \neq j \neq k}\theta_\eta(x_1 - E_1)\theta_\eta(x_2 - E_2)\theta_\eta(x_3 - E_3) = \mathbb{E}_{\mathbf{w}}A_1 + \mathbb{E}_{\mathbf{w}}A_2 + \mathbb{E}_{\mathbf{w}}A_3, \quad (8.7)$$

where

$$\begin{aligned} A_1 &:= \frac{1}{N(N-1)(N-2)}\prod_{j=1}^3\left[\frac{1}{N}\sum_i\theta_\eta(\lambda_i - E_j)\right], \\ A_3 &:= \frac{2}{N(N-1)(N-2)}\sum_i\theta_\eta(\lambda_i - E_1)\theta_\eta(\lambda_i - E_2)\theta_\eta(\lambda_i - E_3) + \dots \end{aligned}$$

and

$$A_2 := B_1 + B_2 + B_3, \quad \text{with} \quad B_3 = -\frac{1}{N(N-1)(N-2)} \sum_i \theta_\eta(\lambda_i - E_1) \theta_\eta(\lambda_i - E_2) \sum_k \theta_\eta(\lambda_k - E_3),$$

and similarly,  $B_1$  consists of terms with  $j = k$ , while  $B_2$  consists of terms with  $i = k$ .

Notice that, modulo a trivial change in the prefactor,  $\mathbb{E}_{\mathbf{w}} A_1$  can be approximated by

$$\mathbb{E}_{\mathbf{w}} F \left( \frac{1}{N} \text{Im Tr} \frac{1}{H^{(v)} - z_1}, \dots, \frac{1}{N} \text{Im Tr} \frac{1}{H^{(v)} - z_3} \right),$$

where the function  $F$  is chosen to be  $F(x_1, x_2, x_3) := x_1 x_2 x_3$  if  $\max_j |x_j| \leq N^\varepsilon$  and it is smoothly cutoff to go to zero in the regime  $\max_j |x_j| \geq N^{2\varepsilon}$ . The difference between the expectation of  $F$  and  $A_1$  is negligible, since it comes from the regime where  $N^\varepsilon \leq \max_j \frac{1}{N} |\text{Im Tr} (H^{(v)} - z_j)^{-1}| \leq N^2$ , which has an exponentially small probability by (8.3) (the upper bound on the Green's function always holds since  $\eta \geq N^{-2}$ ). Here the arguments of  $F$  are imaginary parts of the trace of the Green's function, but this type of function is allowed when applying Theorem 2.3, since

$$\text{Im Tr} G(z) = \frac{1}{2} [\text{Tr} G(z) - \text{Tr} G(\bar{z})].$$

We remark that the main assumption (2.21) for Theorem 2.3 is satisfied by using (2.17) of Theorem 2.1 with the choice of  $M \asymp N$ .

Similarly, we can approximate  $\mathbb{E}_{\mathbf{w}} B_3$  by

$$\mathbb{E}_{\mathbf{w}} G \left( \frac{1}{N^2} \text{Tr} \left\{ \text{Im} \frac{1}{H^{(v)} - z_1} \text{Im} \frac{1}{H^{(v)} - z_2} \right\}, \frac{1}{N} \text{Im Tr} \frac{1}{H^{(v)} - z_3} \right),$$

where  $G(x_1, x_2) = x_1 x_2$  with an appropriate cutoff for large arguments, and there are similar expressions for  $B_1, B_2$  and also for  $A_3$ , the latter involving the trace of the product of three resolvents. By Theorem 2.3, these expectations w.r.t.  $\mathbf{w}$  in the approximations of  $\mathbb{E}_{\mathbf{w}} A_i$  can be replaced by expectations w.r.t.  $\mathbf{v}$  with only negligible errors provided that  $\eta \geq N^{-1-\varepsilon}$ . We have thus proved that

$$\lim_{N \rightarrow \infty} \int dx_1 dx_2 dx_3 [p_{w,N}^{(3)}(x_1, x_2, x_3) - p_{v,N}^{(3)}(x_1, x_2, x_3)] \theta_\eta(x_1 - E_1) \theta_\eta(x_2 - E_2) \theta_\eta(x_3 - E_3) = 0. \quad (8.8)$$

Set  $\eta = N^{-1-\varepsilon}$  for the rest of the proof. We now show that the validity of (8.8) for any choice of  $E, \alpha_1, \alpha_2, \alpha_3$  (recall  $E_j = E + \alpha_j/N$ ) implies that the rescaled correlation functions,  $p_{w,N}^{(3)}(E + \beta_1/N, \dots, E + \beta_3/N)$  and  $p_{v,N}^{(3)}(E + \beta_1/N, \dots, E + \beta_3/N)$ , as functions of the variables  $\beta_1, \beta_2, \beta_3$ , have the same weak limit.

Let  $O$  be a smooth, compactly supported test function and let

$$O_\eta(\beta_1, \beta_2, \beta_3) := \frac{1}{(\pi N)^3} \int_{\mathbb{R}^3} d\alpha_1 d\alpha_2 d\alpha_3 O(\alpha_1, \alpha_2, \alpha_3) \theta_\eta \left( \frac{\beta_1 - \alpha_1}{N} \right) \dots \theta_\eta \left( \frac{\beta_3 - \alpha_3}{N} \right)$$

be its smoothing on scale  $N\eta$ . Then we can write

$$\begin{aligned} & \int_{\mathbb{R}^3} d\beta_1 d\beta_2 d\beta_3 O(\beta_1, \beta_2, \beta_3) p_{w,N}^{(3)} \left( E + \frac{\beta_1}{N}, \dots, E + \frac{\beta_3}{N} \right) \\ &= \int_{\mathbb{R}^3} d\beta_1 d\beta_2 d\beta_3 O_\eta(\beta_1, \beta_2, \beta_3) p_{w,N}^{(3)} \left( E + \frac{\beta_1}{N}, \dots, E + \frac{\beta_3}{N} \right) \\ &+ \int_{\mathbb{R}^3} d\beta_1 d\beta_2 d\beta_3 (O - O_\eta)(\beta_1, \beta_2, \beta_3) p_{w,N}^{(3)} \left( E + \frac{\beta_1}{N}, \dots, E + \frac{\beta_3}{N} \right). \end{aligned} \quad (8.9)$$

The first term on the right side, after the change of variables  $x_j = E + \beta_j/N$ , is equal to

$$\int_{\mathbb{R}^3} d\alpha_1 d\alpha_2 d\alpha_3 O(\alpha_1, \alpha_2, \alpha_3) \int_{\mathbb{R}^3} dx_1 dx_2 dx_3 p_{w,N}^{(3)}(x_1, x_2, x_3) \theta_\eta(x_1 - E_1) \theta_\eta(x_2 - E_2) \theta_\eta(x_3 - E_3), \quad (8.10)$$

i.e., it can be written as an integral of expressions of the form (8.8) for which limits with  $p_{w,N}$  and  $p_{v,N}$  coincide.

Finally, the second term on the right hand side of (8.9) is negligible. To see this, notice that for any test function  $Q$ , we have

$$\begin{aligned} & \int_{\mathbb{R}^3} d\beta_1 d\beta_2 d\beta_3 Q(\beta_1, \beta_2, \beta_3) p_{w,N}^{(3)}\left(E + \frac{\beta_1}{N}, \dots, E + \frac{\beta_3}{N}\right) \\ &= N^3 \int_{\mathbb{R}^3} dx_1 dx_2 dx_3 Q(N(x_1 - E), N(x_2 - E), N(x_3 - E)) p_{w,N}^{(3)}(x_1, x_2, x_3) \\ &= \left(1 - \frac{1}{N}\right) \left(1 - \frac{2}{N}\right) \mathbb{E}_{\mathbf{w}} \sum_{i \neq j \neq k} Q(N(\lambda_i - E), N(\lambda_j - E), N(\lambda_k - E)). \end{aligned} \quad (8.11)$$

If the test function  $Q$  were supported on a ball of size  $N^{\varepsilon'}$ ,  $\varepsilon' > 0$ , then this last term were bounded by

$$\|Q\|_\infty \mathbb{E}_{\mathbf{w}} \mathcal{N}_{CN^{-1+\varepsilon'}}^3(E) \leq C \|Q\|_\infty N^{4\varepsilon'}. \quad (8.12)$$

Here  $\mathcal{N}_\tau(E)$  denotes the number of eigenvalues in the interval  $[E - \tau, E + \tau]$  and in the estimate we used the local semicircle law on intervals of size  $\tau \geq N^{-1+\varepsilon'}$ .

Set now  $Q := O - O_\eta$ . From the definition of  $O_\eta$ , it is easy to see that the function

$$Q_1(\beta_1, \beta_2, \beta_3) = O(\beta_1, \beta_2, \beta_3) - O_\eta(\beta_1, \beta_2, \beta_3) \prod_{j=1}^3 1(|\beta_j| \leq N^{\varepsilon'})$$

satisfies the bound  $\|Q_1\|_\infty \leq \|Q\|_\infty = \|O - O_\eta\|_\infty \leq CN\eta = CN^{-\varepsilon}$ . So choosing  $\varepsilon' < \varepsilon/4$ , the contribution of  $Q_1$  is negligible. Finally,  $Q_2 = Q - Q_1$  is given by

$$Q_2(\beta_1, \beta_2, \beta_3) = -O_\eta(\beta_1, \beta_2, \beta_3) \left[ 1 - \prod_{j=1}^3 1(|\beta_j| \leq N^{\varepsilon'}) \right]$$

and

$$\begin{aligned} |Q_2| &\leq C \left[ \frac{1}{1 + \beta_1^2} \right] \left[ \frac{1}{1 + \beta_2^2} \right] \left[ \frac{1}{1 + \beta_3^2} \right] \left\{ 1(|\beta_1| \geq N^{\varepsilon'}) + \dots \right\} \\ &\leq C \left\{ N^{-\varepsilon'} \left[ \frac{N^{\varepsilon'}}{N^{2\varepsilon'} + \beta_1^2} \right] \left[ \frac{1}{1 + \beta_2^2} \right] \left[ \frac{1}{1 + \beta_3^2} \right] + \dots \right\}. \end{aligned} \quad (8.13)$$

Hence the contribution of  $Q_2$  in the last term of (8.11) is bounded by

$$CN^{-3-\varepsilon'} \mathbb{E}_{\mathbf{w}} \sum_{i,j,k} \left\{ \left[ \frac{N^{-1+\varepsilon'}}{N^{-2+2\varepsilon'} + (\lambda_i - E)^2} \right] \left[ \frac{N^{-1}}{N^{-2} + (\lambda_j - E)^2} \right] \left[ \frac{N^{-1}}{N^{-2} + (\lambda_k - E)^2} \right] + \dots \right\}$$

From Theorem 2.1, the last term is bounded by  $N^{-\varepsilon'}$  up to some logarithmic factor. This completes the proof of Theorem 6.4.  $\square$



## A Spectral condition for band matrices

**Lemma A.1** *Let  $B = (\sigma_{ij})$  satisfying (2.1) and (2.8) with  $W \geq 1$  and with  $f$  being a nonnegative symmetric function with  $\int f = 1$  and  $f \in L^\infty(\mathbb{R})$ . Then we have*

$$B \geq -1 + \delta \tag{A.1}$$

for some  $\delta > 0$  and  $W$  large enough, depending on  $f$ .

*Proof.* Recall that the discrete Fourier transform in  $d = 1$  dimensions is defined as follows. Let  $\varepsilon := 1/N$  and

$$\Lambda_\varepsilon := \Lambda = \varepsilon\mathbb{Z}/\mathbb{Z}$$

be the periodic one dimensional lattice (torus) of size 1 and spacing  $\varepsilon$  with its dual lattice being

$$\Lambda_\varepsilon^* := \Lambda^* := \left(2\pi\mathbb{Z}\right) / \left(\frac{2\pi}{\varepsilon}\mathbb{Z}\right).$$

Let  $\psi$  be a function on  $\Lambda$ . Then its Fourier transform  $\mathcal{F}_N\psi$  is a function on  $\Lambda^*$  defined as

$$\mathcal{F}_N\psi(p) = \varepsilon \sum_{x \in \Lambda} \psi(x) e^{-ip \cdot x}$$

and it is an isometry

$$\varepsilon \sum_{x \in \Lambda} \overline{\psi(x)} \phi(x) = \sum_{p \in \Lambda^*} \overline{\mathcal{F}_N\psi(p)} \mathcal{F}_N\phi(p).$$

In our case,  $x = k/N$  and define

$$F_W(x) := NW^{-1}f(xN/W).$$

Then, for  $p \in \Lambda^*$ , we have

$$(\mathcal{F}_N F_W)(p) = \sum_{x \in \Lambda} F_W(x) e^{-ip \cdot x} = \sum_{k=1}^N W^{-1}f(k/W) e^{-i(Wp/N) \cdot (k/W)} = \widehat{f}(q) + o(1), \quad q = Wp/N,$$

where the error term vanishes as  $W \rightarrow \infty$  and  $\widehat{f}$  denotes the usual Fourier transform in  $L^1(\mathbb{R})$

$$\widehat{f}(q) = \int f(y) e^{-iqy} dy.$$

With this formula, and with the notation  $\psi_N(j) := \psi(j/N)$  for any  $\psi$  defined on  $\Lambda$ , we have

$$(B\psi_N)(k) = \varepsilon \sum_{\ell=1}^N NW^{-1}f((k-\ell)/W)\psi_N(\ell) = \varepsilon \sum_{y \in \Lambda} F_W\left(\frac{k}{N} - y\right)\psi(y) = \sum_{p \in \Lambda^*} e^{ipk/N} \mathcal{F}_N F_W(p) \mathcal{F}_N\psi(p)$$

and

$$\sum_{p \in \Lambda^*} |\mathcal{F}_N\psi(p)|^2 = \varepsilon \sum_{j=1}^N |\psi_N(j)|^2$$

which is normalized to be 1. Hence  $B$  on the Fourier side acts as a multiplication by the function  $\mathcal{F}_N F_W$ , so

$$\text{Spec}B = \text{Range } \mathcal{F}_N F_W \subset \text{supp } \widehat{f} + o(1).$$

Since  $f$  is nonnegative, symmetric function and  $\int f = 1$ , we have  $\widehat{f}$  is real and

$$\inf \widehat{f} > -1 + \delta$$

for some  $\delta > 0$ , which completes the proof.

## B Large deviation estimates

In this Appendix we prove two large deviations results. They are weaker than the corresponding results of Hanson and Wright [22], used in [14], but they require only independent, not necessarily identically distributed random variables, moreover the proofs are much simpler.

**Lemma B.1** *Let  $a_i$  ( $1 \leq i \leq N$ ) be  $N$  independent complex random variables with mean zero, variance  $\sigma^2$  and uniform subexponential decay, i.e., there exist  $\alpha, \beta > 0$  that for any  $x > 0$*

$$\mathbb{P}(|a_i| \geq x^\alpha) \leq \beta e^{-x}. \quad (\text{B.1})$$

Then for any  $A_i \in \mathbb{C}$  ( $1 \leq i \leq N$ ) and  $D \geq 1$  we have,

$$\mathbb{P} \left\{ \left| \sum_i a_i A_i \right| \geq D\sigma \left( \sum_i |A_i|^2 \right)^{1/2} \right\} \leq C \exp(-cD^{\frac{2}{2+\alpha}}) \quad (\text{B.2})$$

for some positive constants  $C$  and  $c$  depending on  $\alpha$  and  $\beta$  in (B.1).

*Proof of Lemma B.1.* Without loss of generality, we may assume that  $\sigma = 1$ . The assumption (B.1) implies that the  $k$ -th moment of  $a_i$  is bounded by:

$$\mathbb{E}|a_i|^k \leq (Ck)^{\alpha k} \quad (\text{B.3})$$

for some  $C > 0$  depending on  $\alpha$  and  $\beta$ .

First, for  $p \in \mathbb{N}$ , we estimate

$$\mathbb{E} \left| \sum_{i=1}^N a_i A_i \right|^p. \quad (\text{B.4})$$

With the Marcinkiewicz-Zygmund inequality, for an integer  $p \geq 2$ , we have

$$\mathbb{E} \left| \sum_i a_i A_i \right|^p \leq (Cp)^{p/2} \mathbb{E} \left[ \left( \sum_i |a_i A_i|^2 \right)^{p/2} \right] \quad (\text{B.5})$$

(for the estimate of the constant, see e.g. Exercise 2.2.30 of [30]). Using (B.3), we have  $\mathbb{E}|a_{i_1} a_{i_2} \cdots a_{i_{p/2}}|^2 \leq (Cp)^{\alpha p}$ . Inserting it into (B.5), we obtain

$$\mathbb{E} \left| \sum_i a_i A_i \right|^p \leq (Cp^{\frac{1}{2}+\alpha})^p \left( \sum_i |A_i|^2 \right)^{p/2}, \quad (\text{B.6})$$

which implies (B.2) by choosing an even integer  $p$  of the order  $(D/Ce)^{\frac{2}{2+\alpha}}$  and applying a high moment Markov inequality.  $\square$

**Lemma B.2** *Let  $a_i$  ( $1 \leq i \leq N$ ) be  $N$  independent random complex variables with mean zero, variance  $\sigma^2$  and having the uniform subexponential decay (B.1). Let  $B_{ij} \in \mathbb{C}$  ( $1 \leq i, j \leq N$ ). Then we have that*

$$\mathbb{P} \left\{ \left| \sum_{i=1}^N \bar{a}_i B_{ii} a_i - \sum_{i=1}^N \sigma^2 B_{ii} \right| \geq D\sigma^2 \left( \sum_{i=1}^N |B_{ii}|^2 \right)^{1/2} \right\} \leq C \exp(-cD^{\frac{1}{1+\alpha}}) \quad (\text{B.7})$$

and

$$\mathbb{P} \left\{ \left| \sum_{i \neq j} \bar{a}_i B_{ij} a_j \right| \geq D\sigma^2 \left( \sum_{i \neq j} |B_{ij}|^2 \right)^{1/2} \right\} \leq C \exp(-cD^{\frac{1}{2(1+\alpha)}}) \quad (\text{B.8})$$

for some positive constants  $C$  and  $c$  depending on  $\alpha$  and  $\beta$  in (B.1).

*Proof of Lemma B.2.* Without loss of generality, we may again assume that  $\sigma = 1$ . First, we prove (B.7). Notice that  $|a_i|^2 - 1$  ( $1 \leq i \leq N$ ) are independent random variables with mean 0 and variance less than some constants  $C$ . Furthermore, the  $k$ -th moment of  $|a_i|^2 - 1$  is bounded as

$$\mathbb{E}(|a_i|^2 - 1)^k \leq (Ck)^{2\alpha k}. \quad (\text{B.9})$$

Then following the proof of the Lemma B.1 with  $|a_i|^2 - 1$  replacing  $a_i$ , we obtain (B.7).

Next, we prove (B.8). For any  $p \in \mathbb{N}$ ,  $p \geq 2$ , we estimate

$$\mathbb{E} \left| \sum_i \bar{a}_i \xi_i \right|^p \equiv \mathbb{E} \left| \sum_{i>j} \bar{a}_i B_{ij} a_j \right|^p \quad (\text{B.10})$$

where  $\xi_i := \sum_{j<i} B_{ij} a_j$ . Note that  $a_i$  and  $\xi_i$  are independent for any fixed  $i$ . By the definition,

$$X_n \equiv \sum_{i=1}^n \bar{a}_i \xi_i \quad (\text{B.11})$$

is martingale. Using the Burkholder inequality, we have that

$$\mathbb{E} \left| \sum_i \bar{a}_i \xi_i \right|^p \leq (Cp)^{3p/2} \mathbb{E} \left[ \left( \sum_i |\bar{a}_i \xi_i|^2 \right)^{p/2} \right] \quad (\text{B.12})$$

(for the constant, see Section VII.3 of [28]). By the generalized Minkowski inequality, by the independence of  $a_i$  and  $\xi_i$  and using (B.3), we have

$$\left[ \mathbb{E} \left( \sum_i |\bar{a}_i \xi_i|^2 \right)^{p/2} \right]^{2/p} \leq \sum_i \left[ \mathbb{E} |\bar{a}_i \xi_i|^p \right]^{2/p} = \sum_i \left[ \mathbb{E} (|\bar{a}_i|^p) \mathbb{E} (|\xi_i|^p) \right]^{2/p} \leq (Cp)^{2\alpha} \sum_i \left[ \mathbb{E} (|\xi_i|^p) \right]^{2/p}.$$

Using (B.6), we have

$$\mathbb{E}(|\xi_i|^p) \leq (Cp^{\frac{1}{2}+\alpha})^p \left( \sum_j |B_{ij}|^2 \right)^{p/2}.$$

Combining this with (B.12) we obtain

$$\mathbb{E} \left| \sum_i \bar{a}_i \xi_i \right|^p \leq (Cp)^{2p(1+\alpha)} \left( \sum_i \sum_j |B_{ij}|^2 \right)^{p/2}. \quad (\text{B.13})$$

Then choosing  $(D/Ce)^{\frac{1}{2(1+\alpha)}}$  and applying Markov inequality, we obtain (B.8).  $\square$

In our applications we will need these two lemmas when  $D$  is a power of  $\log N$ . For simplicity, we do not want to keep track of the precise powers in the estimate and we are interested only in error bounds that decay faster than any fixed power of  $N$ , say  $CN^{-\log \log N}$ . Therefore, in this paper we will use the following weaker form of these two lemmas, the stronger form will be useful in future applications.

**Corollary B.3** *Let  $a_i$  ( $1 \leq i \leq N$ ) be  $N$  independent random complex variables with mean zero, variance  $\sigma^2$  and having the uniform subexponential decay (B.1). Let  $A_i, B_{ij} \in \mathbb{C}$  ( $1 \leq i, j \leq N$ ). Then we have that*

$$\mathbb{P} \left\{ \left| \sum_{i=1}^N a_i A_i \right| \geq (\log N)^{\frac{3}{2}+\alpha} \sigma \left( \sum_i |A_i|^2 \right)^{1/2} \right\} \leq CN^{-\log \log N}, \quad (\text{B.14})$$

$$\mathbb{P} \left\{ \left| \sum_{i=1}^N \bar{a}_i B_{ii} a_i - \sum_{i=1}^N \sigma^2 B_{ii} \right| \geq (\log N)^{\frac{3}{2}+2\alpha} \sigma^2 \left( \sum_{i=1}^N |B_{ii}|^2 \right)^{1/2} \right\} \leq CN^{-\log \log N}, \quad (\text{B.15})$$

$$\mathbb{P} \left\{ \left| \sum_{i \neq j} \bar{a}_i B_{ij} a_j \right| \geq (\log N)^{3+2\alpha} \sigma^2 \left( \sum_{i \neq j} |B_{ij}|^2 \right)^{1/2} \right\} \leq CN^{-\log \log N}, \quad (\text{B.16})$$

for some constants  $C$  depending on  $\alpha$  and  $\beta$  in (B.1).

## C Proof of Lemma 6.5

We first prove a version of this lemma when the fourth moment exactly matches, i.e.,  $\gamma = 0$ , then we explain how to deal with the approximation. More precisely, we first show the following:

**Lemma C.1** *Under the condition*

$$m_4 - m_3^2 - 1 \geq C_1, \quad m_4 \leq C_2 \quad (\text{C.1})$$

for some positive constants  $C_1$  and  $C_2$ , there exists a real random variable  $\xi$  such that the first four moments of  $\xi$  are 0, 1,  $m_3$  and  $m_4$  and the distribution  $\nu$  of  $\xi$  satisfies logarithmic Sobolev inequality and the LSI constant is bounded from above by a function of  $C_1$  and  $C_2$ . Moreover,  $\nu$  can be chosen to be absolutely continuous with a smooth positive density,  $\nu(dx) = e^{-U(x)} dx$ , such that the derivatives of  $U$  satisfy

$$|U^{(k)}(x)| \leq C_k (1+x^2)^{C_k} \quad (\text{C.2})$$

with some fixed constant  $C$  and  $k$ -dependent constants  $C_k$ .

*Remark.* The last statement about the smoothness of  $U$  will not be needed in this paper, but we state it for further reference.

*Proof.* We start with the case  $|m_3| > \delta$ , where  $\delta$  is small enough number to depend only on  $C_1$ , see below. Let  $\xi$  be the sum of two Gaussians, with density function of the form

$$f_\xi(x) = \frac{b}{(a+b)} \frac{1}{\sqrt{2\pi\sigma}} e^{-(x-a)^2/(2\sigma)} + \frac{a}{(a+b)} \frac{1}{\sqrt{2\pi\sigma}} e^{-(x+b)^2/(2\sigma)} \quad (\text{C.3})$$

with some parameters  $a > 0$ ,  $b > 0$ ,  $\sigma > 0$ . If the first 4 moments of  $f_\xi(x)$  are 0, 1,  $m_3$  and  $m_4$ , then we have the relations

$$m_4 = 1 + \frac{m_3^2}{1-\sigma} + 4\sigma - 2\sigma^2, \quad (\text{C.4})$$

$$ab = 1 - \sigma \quad \text{and} \quad a - b = \frac{m_3}{1-\sigma}. \quad (\text{C.5})$$

With  $m_3, m_4$  in (C.1) and  $|m_3| \geq \delta$ , one can always find a solution of (C.4) such that  $0 < \sigma < 1$ . Actually, one can see that  $c < \sigma < C$ , where  $c$  and  $C$  only depend on  $C_1, C_2$  in (6.19) and  $\delta$ .

Once  $\sigma$  is found, it is easy to check that one can always find real solutions  $a, b$  for (C.5) as long as  $m_3, m_4$  satisfy (C.1) and  $|m_3| \geq \delta$ . Since the solutions  $a, b, \sigma$  are continuous with respect to  $m_3$  and  $m_4$ , then they are uniformly bounded. Distributions of the form (C.3) satisfy the LSI, since they are log concave away from a compact set. Since the parameters  $a, b, \sigma$  are in a compact set, the LSI constant will remain uniformly bounded with a bound depending on  $C_1, C_2$  and  $\delta$ . It is clear that the density function (C.3) is positive and its logarithm satisfies (C.2).

Now we consider the case that  $|m_3| < \delta$  with a small  $\delta = \frac{1}{100} \min\{1, C_1\}$ , where  $C_1$  is the constant in (C.1). Without loss of generality, we may assume  $m_3 > 0$ . We consider the following three parameter family of probability densities

$$f_{d,\beta,\varepsilon}(x) = (1-\varepsilon)g_{d,\beta}(x) + \varepsilon h(x)$$

with

$$g_{d,\beta}(x) = \frac{\beta+1}{2d^{\beta+1}} \cdot |x|^\beta \cdot \mathbf{1}(|x| \leq d), \quad h(x) = \frac{b}{(a+b)} \frac{1}{\sqrt{2\pi}} e^{-(x-a)^2/2} + \frac{a}{(a+b)} \frac{1}{\sqrt{2\pi}} e^{-(x+b)^2/2},$$

where the parameters are in the range  $-1 < \beta < \infty$ ,  $0 < d < \infty$ ,  $0 \leq \varepsilon \ll 1$  and  $a, b$  will be chosen explicitly. Simple calculation shows that the moments of  $f_{d,\beta,\varepsilon}$  are  $m_1 = 0$ ,

$$m_2 = (1-\varepsilon) \frac{\beta+1}{\beta+3} d^2 + \varepsilon(1+ab), \quad (\text{C.6})$$

$$m_3 = \varepsilon ab(a-b), \quad (\text{C.7})$$

$$m_4 = (1-\varepsilon) \frac{\beta+1}{\beta+5} d^4 + \varepsilon \left[ 3 + ab(6 + a^2 + b^2 - ab) \right]. \quad (\text{C.8})$$

Choosing, say,  $a = 2$ ,  $b = 1$ , and setting  $m_2 = 1$ , we obtain  $d^2 = \frac{1-3\varepsilon}{1-\varepsilon} \frac{\beta+3}{\beta+1}$  from the first equation,  $\varepsilon = m_3/2$  from the second equation and finally the last equation becomes

$$m_4 = \frac{(1-3m_3/2)^2}{1-m_3/2} \frac{(\beta+3)^2}{(\beta+1)(\beta+5)} + \frac{23}{2} m_3. \quad (\text{C.9})$$

Recall that we are in the regime where  $|m_3| \leq \delta \leq C_1/100$ . For any fixed  $0 \leq m_3 \leq \delta$ , the right hand side of (C.9) is a monotonically decreasing function in  $\beta \in (-1, \infty)$  whose value goes down from  $\infty$  to  $\frac{(1-3m_3/2)^2}{1-m_3/2} + \frac{23}{2}m_3 \leq 1 + 20\delta$ . But we know from (C.1) that  $C_2 \geq m_4 \geq 1 + 100\delta$ , thus there is a value  $\beta$  such that (C.9) holds, moreover,  $\beta$  is in a compact subinterval of  $(-1, \infty)$  that depends only on  $\delta$  and  $C_2$ . It is then easy to check that the support and the supremum norm of the density  $g_{d,\beta}$  also remains in a compact set, depending only on  $\delta$ . Therefore we constructed a probability measure with the given moments, that is a linear combination of two Gaussians plus a compactly supported piece with a nonnegative bounded density. To ensure smoothness, we replace  $g_{d,\beta}$  with  $\tilde{g}_{d,\beta,\tau} := \vartheta_\tau * g_{d,\beta}$ , where  $\vartheta_\tau(x) = \tau^{-1}\vartheta(x/\tau)$  and  $\vartheta$  is a compactly supported nonnegative smooth symmetric function with  $\int \vartheta = 1$ . The first moment  $m_1$  is unchanged and the formulas (C.6) for the higher moments will get modified by an error term of order  $\tau$ . Let  $\tau$  be much smaller than all other parameters in this proof. It is easy to see that, by a simple calculation treating  $\tau$  as a small perturbation, one can still choose  $a, b, \varepsilon$  and  $\beta$  in the previous argument to match  $m_2 = 1, m_3$  and  $m_4$ .

Finally, note that the sum of two Gaussians satisfy the LSI, as well as its compact perturbation and the new LSI constant depends only on the supremum norm of the density of the perturbation. Since all these parameters remain uniformly controlled by  $C_1$  and  $C_2$ , we proved Lemma C.1, i.e., Lemma 6.5 for  $\gamma = 0$ .  $\square$

Now consider the case  $\gamma > 0$ . For any real random variable  $\zeta$ , independent of  $\xi^G$ , and with the first 4 moments being 0, 1,  $m_3(\zeta)$  and  $m_4(\zeta) < \infty$ , the first 4 moments of

$$\zeta' = (1 - \gamma)^{1/2}\zeta + \gamma^{1/2}\xi^G \quad (\text{C.10})$$

are 0, 1,

$$m_3(\zeta') = (1 - \gamma)^{3/2}m_3(\zeta) \quad (\text{C.11})$$

and

$$m_4(\zeta') = (1 - \gamma)^2m_4(\zeta) + 6\gamma - 3\gamma^2. \quad (\text{C.12})$$

Given  $m_3$  and  $m_4$ , satisfying (C.1) and using Lemma C.1, we obtain that for any  $\gamma$  small enough, there exists a real random variable  $\xi_\gamma$  such that the first four moments are 0, 1,

$$m_3(\xi_\gamma) = (1 - \gamma)^{-3/2}m_3 \quad (\text{C.13})$$

and

$$m_4(\xi_\gamma) = m_3(\xi_\gamma)^2 + (m_4 - m_3^2).$$

With  $m_4 \leq C_2$ , we have  $m_3^2 \leq C_2$ , thus

$$|m_4(\xi_\gamma) - m_4| \leq C\gamma \quad (\text{C.14})$$

for some  $C$  depending on  $C_2$ .

Hence with (C.11) and (C.12), we obtain that  $\xi' = (1 - \gamma)^{1/2}\xi_\gamma + \gamma^{1/2}\xi^G$  satisfies  $m_3(\xi') = m_3$  and (6.21). With Lemma C.1, we obtain that the LSI constant of  $\xi_\gamma$  is bounded by a constant only depends on  $C_1$  and  $C_2$ , which completes the proof of Lemma 6.5.  $\square$

## References

- [1] Anderson, G.; Zeitouni, O. : A CLT for a band matrix model. *Probab. Theory Related Fields* **134** (2006), no. 2, 283–338.
- [2] Ben Arous, G., Péché, S.: Universality of local eigenvalue statistics for some sample covariance matrices. *Comm. Pure Appl. Math.* **LVIII**. (2005), 1–42.
- [3] Bleher, P., Its, A.: Semiclassical asymptotics of orthogonal polynomials, Riemann-Hilbert problem, and universality in the matrix model. *Ann. of Math.* **150** (1999): 185–266.
- [4] Bobkov, S. G., Götze, F.: Exponential integrability and transportation cost related to logarithmic Sobolev inequalities. *J. Funct. Anal.* **163** (1999), no. 1, 1–28.
- [5] Brézin, E., Hikami, S.: Correlations of nearby levels induced by a random potential. *Nucl. Phys. B* **479** (1996), 697–706, and Spectral form factor in a random matrix theory. *Phys. Rev. E* **55** (1997), 4067–4083.
- [6] Deift, P.: Orthogonal polynomials and random matrices: a Riemann-Hilbert approach. *Courant Lecture Notes in Mathematics* **3**, American Mathematical Society, Providence, RI, 1999
- [7] Deift, P., Gioev, D.: Random Matrix Theory: Invariant Ensembles and Universality. *Courant Lecture Notes in Mathematics* **18**, American Mathematical Society, Providence, RI, 2009
- [8] Deift, P., Kriecherbauer, T., McLaughlin, K.T-R, Venakides, S., Zhou, X.: Uniform asymptotics for polynomials orthogonal with respect to varying exponential weights and applications to universality questions in random matrix theory. *Comm. Pure Appl. Math.* **52** (1999):1335–1425.
- [9] Deift, P., Kriecherbauer, T., McLaughlin, K.T-R, Venakides, S., Zhou, X.: Strong asymptotics of orthogonal polynomials with respect to exponential weights. *Comm. Pure Appl. Math.* **52** (1999): 1491–1552.
- [10] Disertori, M., Pinson, H., Spencer, T.: Density of states for random band matrices. *Commun. Math. Phys.* **232**, 83–124 (2002)
- [11] Dyson, F.J.: Correlations between eigenvalues of a random matrix. *Commun. Math. Phys.* **19**, 235-250 (1970).
- [12] Erdős, L., Schlein, B., Yau, H.-T.: Semicircle law on short scales and delocalization of eigenvectors for Wigner random matrices. *Ann. Probab.* **37**, No. 3, 815–852 (2008)
- [13] Erdős, L., Schlein, B., Yau, H.-T.: Local semicircle law and complete delocalization for Wigner random matrices. *Commun. Math. Phys.* **287**, 641–655 (2009)
- [14] Erdős, L., Schlein, B., Yau, H.-T.: Wegner estimate and level repulsion for Wigner random matrices. *Int Math Res Notices* **2010** (3): 436-479 (2010)
- [15] Erdős, L., Schlein, B., Yau, H.-T.: Universality of random matrices and local relaxation flow. To appear in *Invent. Math.* [arxiv.org/abs/0907.5605](https://arxiv.org/abs/0907.5605)
- [16] Erdős, L., Ramirez, J., Schlein, B., Yau, H.-T.: *Universality of sine-kernel for Wigner matrices with a small Gaussian perturbation.* *Electr. J. Prob.* **15**, Paper 18, 526–604 (2010)

- [17] Erdős, L., Péché, G., Ramírez, J., Schlein, B., and Yau, H.-T., Bulk universality for Wigner matrices. *Comm. Pure Appl. Math.* **63**, No. 7, 895-925 (2010)
- [18] Erdős, L., Ramírez, J., Schlein, B., Tao, T., Vu, V. and Yau, H.-T., Bulk universality for Wigner hermitian matrices with subexponential decay. *Math. Res. Lett.* **17** (4), 667–674 (2010).
- [19] Erdős, L., Schlein, B., Yau, H.-T., Yin, J.: The local relaxation flow approach to universality of the local statistics for random matrices. To appear in *Annales Inst. H. Poincaré, Prob. and Stat.* Preprint arXiv:0911.3687
- [20] Erdős, L., Yau, H.-T., Yin, J.: Universality for generalized Wigner matrices with Bernoulli distribution. To appear in *Journal of Combinatorics*. Preprint arXiv:1003.3813
- [21] Guionnet, A.: Large deviation upper bounds and central limit theorems for band matrices, *Ann. Inst. H. Poincaré Probab. Statist* **38** , (2002), pp. 341-384.
- [22] Hanson, D.L., Wright, F.T.: A bound on tail probabilities for quadratic forms in independent random variables. *The Annals of Math. Stat.* **42** (1971), no.3, 1079-1083.
- [23] Johansson, K.: Universality of the local spacing distribution in certain ensembles of Hermitian Wigner matrices. *Comm. Math. Phys.* **215** (2001), no.3. 683–705.
- [24] Johansson, K.: Universality for certain hermitian Wigner matrices under weak moment conditions. Preprint arxiv.org/abs/0910.4467
- [25] Mehta, M.L.: *Random Matrices*. Academic Press, New York, 1991.
- [26] Mehta, M.L., Gaudin, M.: On the density of eigenvalues of a random matrix. *Nuclear Phys.* **18**, 420-427 (1960).
- [27] Pastur, L., Shcherbina M.: Bulk universality and related properties of Hermitian matrix models. *J. Stat. Phys.* **130** (2008), no.2., 205-250.
- [28] Shiryaev, A. N.: *Probability*. Graduate Text in Mathematics. **54**. Springer, 1984.
- [29] Spencer, T.: *Random banded and sparse matrices (Chapter 23)* to appear in “Oxford Handbook of Random Matrix Theory”, edited by G. Akemann, J. Baik and P. Di Francesco.
- [30] Stroock, D.W.: *Probability theory, an analytic view*. Cambridge University Press, 1993.
- [31] Tao, T. and Vu, V.: Random matrices: Universality of the local eigenvalue statistics. Preprint arXiv:0906.0510.
- [32] Tao, T. and Vu, V.: Random matrices: Universality of local eigenvalue statistics up to the edge. Preprint. arXiv:0908.1982
- [33] Tao, T. and Vu, V.: Random covariance matrices: Universality of local statistics of eigenvalues. Preprint. arXiv:0912.0966
- [34] Vu, V.: Spectral norm of random matrices. *Combinatorica*, **27** (6) (2007), 721-736.