

# Local circular law for the product of a deterministic matrix with a random matrix

Haokai Xi <sup>\*1</sup>, Fan Yang <sup>†1</sup> and Jun Yin <sup>‡1</sup>

<sup>1</sup>Department of Mathematics, University of Wisconsin-Madison

June 7, 2016

## Abstract

It is well known that the spectral measure of eigenvalues of a rescaled square non-Hermitian random matrix with independent entries satisfies the circular law. We consider the product  $TX$ , where  $T$  is a deterministic  $N \times M$  matrix and  $X$  is a random  $M \times N$  matrix with independent entries having zero mean and variance  $(N \wedge M)^{-1}$ . We prove a general local circular law for the empirical spectral distribution (ESD) of  $TX$  at any point  $z$  away from the unit circle under the assumptions that  $N \sim M$ , and the matrix entries  $X_{ij}$  have sufficiently high moments. More precisely, if  $z$  satisfies  $||z| - 1| \geq \tau$  for arbitrarily small  $\tau > 0$ , the ESD of  $TX$  converges to  $\tilde{\chi}_{\mathbb{D}}(z)dA(z)$ , where  $\tilde{\chi}_{\mathbb{D}}$  is a rotation-invariant function determined by the singular values of  $T$  and  $dA$  denotes the Lebesgue measure on  $\mathbb{C}$ . The local circular law is valid around  $z$  up to scale  $(N \wedge M)^{-1/4+\epsilon}$  for any  $\epsilon > 0$ . Moreover, if  $|z| > 1$  or the matrix entries of  $X$  have vanishing third moments, the local circular law is valid around  $z$  up to scale  $(N \wedge M)^{-1/2+\epsilon}$  for any  $\epsilon > 0$ .

## 1 Introduction

**Circular law for non-Hermitian random matrices.** The study of the eigenvalue spectral of non-Hermitian random matrices goes back to the celebrated paper [19] by Ginibre, where he calculated the joint probability density for the eigenvalues of non-Hermitian random matrix with independent complex Gaussian entries. The joint density distribution is integrable with an explicit kernel (see [19, 28]), which allowed him to derive the circular law for the eigenvalues. For the Gaussian random matrix with real entries, the joint distribution of the eigenvalues is more complicated but still integrable, which leads to a proof of the circular law as well [6, 10, 18, 35].

For the random matrix with non-Gaussian entries, there is no explicit formula for the joint distribution of the eigenvalues. However, in many cases the eigenvalue spectrum of the non-Gaussian random matrices behaves similarly to the Gaussian case as  $N \rightarrow \infty$ , known as the universality phenomena. A key step in this direction is made by Girko in [20], where he partially proved the circular law for non-Hermitian matrices with independent entries. The crucial insight of the paper is the *Hermitization technique*, which allowed Girko to translate the convergence of complex empirical

---

\*E-mail: haokai@math.wisc.edu.

†E-mail: fyang75@math.wisc.edu.

‡E-mail: jyin@math.wisc.edu. Partially supported by NSF Career Grant DMS-1552192 and Sloan fellowship.

measures of a non-Hermitian matrix into the convergence of logarithmic transforms for a family of Hermitian matrices, or, to be more precise,

$$\mathrm{Tr} \log[(X - z)^\dagger(X - z)] = \log [\det((X - z)^\dagger(X - z))], \quad (1.1)$$

with  $X$  being the random matrix and  $z \in \mathbb{C}$ . Due to the singularity of the log function at 0, the small eigenvalues of  $(X - z)^\dagger(X - z)$  play a special role. The estimate on the smallest singular value of  $X - z$  was not obtained in [20], but the gap was remedied later in a series of paper. Bai [1, 2] analyzed the ESD of  $(X - z)^\dagger(X - z)$  through its Stieltjes transform and handled the logarithmic singularity by assuming bounded density and bounded high moments for the entries of  $X$ . Lower bounds on the smallest singular values were given by Rudelson and Vershynin [31, 32], and subsequently by Tao and Vu [36], Pan and Zhou [30] and Götze and Tikhomirov [21] under weakened moments and smoothness assumptions. The final result was presented in [38], where the circular law is proved under the optimal  $L^2$  assumption. These papers studied the circular law in the global regime, i.e. the convergence of ESD on subsets containing  $\eta N$  eigenvalues for some small constant  $\eta > 0$ . Later in a series of papers [7, 8, 39], Bourgade, Yau and Yin proved the *local* version of the circular law up to the optimal scale  $N^{-1/2+\epsilon}$  under the assumption that the distributions of the matrix entries satisfy a uniform sub-exponential decay condition. In [37], the local universality was proved by Tao and Vu under the assumption of first four moments matching the moments of a Gaussian random variable.

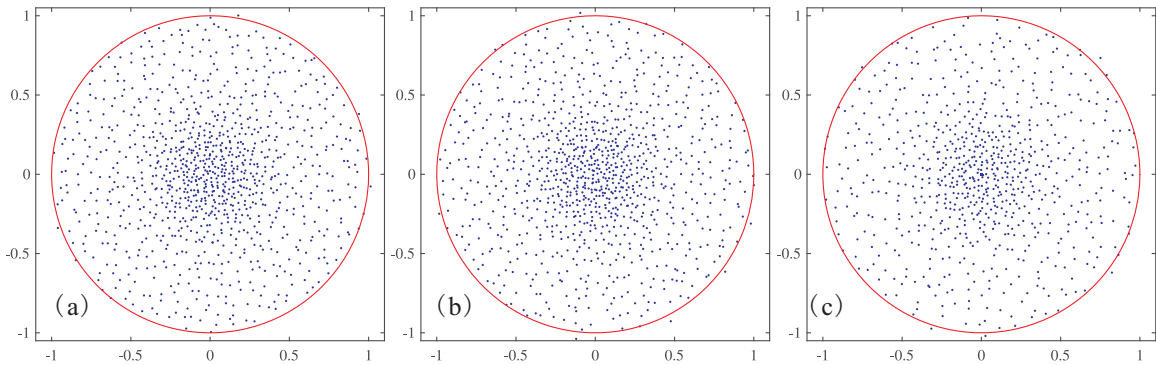


Figure 1: The eigenvalue distribution of the product  $TX$  of a deterministic  $N \times M$  matrix  $T$  with a Gaussian random  $M \times N$  matrix  $X$ . The entries of  $X$  have zero mean and variance  $(N \wedge M)^{-1}$ , and  $TT^\dagger$  has  $0.5(N \wedge M)$  eigenvalues as  $2/17$  and  $0.5(N \wedge M)$  eigenvalues as  $32/17$ . (a)  $N = M = 1000$ . (b)  $N = 1000$ ,  $M = 2000$ . (c)  $N = 1500$ ,  $M = 750$ .

In this paper, we study the ESD of the product of a deterministic  $N \times M$  matrix  $T$  with a random  $M \times N$  matrix  $X$ , where we assume  $N \sim M$ . In Figure 1, we plot the eigenvalue distribution of  $TX$  when  $T$  have two distinct singular values (except the trivial zero singular values). The goal of this paper is to prove a local circular law for the ESD of  $TX$  at any point  $z$  away from the unit circle. Following the idea in [7], the key ingredients for the proof are (a) the upper bound for the largest singular value of  $TX - z$ , (b) the lower bound for the least singular value of  $TX - z$ , and (c) rigidity of the singular values of  $TX - z$ . The upper bound for the largest singular value can be obtained by controlling the norm of  $TX - z$  through a standard large deviation estimate (see e.g. [9, 27, 33] and (2.64)). The lower bound for the least singular value of  $TX - z$  follows from the results in e.g. [32] and [36] (see also Lemma 2.23). Thus the bulk of this paper is devoted to establish (c).

**Basic ideas.** To obtain the rigidity of the singular values of  $TX - z$ , we study the ESD of  $Q := (TX - z)^\dagger(TX - z)$  using Stieltjes transform as in [7]. We normalize  $X$  so that its entries have variance  $(N \wedge M)^{-1}$ . Then  $Q$  is an  $N \times N$  Hermitian matrix with eigenvalues being typically of order 1. We denote its resolvent by  $R(w) := (Q - w)^{-1}$ , where  $w = E + i\eta$  is a spectral parameter with positive imaginary part  $\eta$ . Then the Stieltjes transform of the ESD of  $Q$  is equal to  $N^{-1}\text{Tr } R(w)$ , and we have the convergence estimate

$$N^{-1}\text{Tr } R(w) \approx m_c(w) \quad (1.2)$$

with high probability for large  $N$ . Here  $m_c$  is the Stieltjes transform of the asymptotic eigenvalue density, and the convergence in (1.2) is referred to as the *averaged law*. By taking the imaginary part of (1.2), it is easy to see that a control of the Stieltjes transform yields a control of the eigenvalue density on a small scale of order  $\eta$  around  $E$  (which contains an order  $\eta N$  eigenvalues). A *local law* is an estimate of the form (1.2) for all  $\eta \gg N^{-1}$ . Such local laws have been a cornerstone of the modern random matrix theory. In [16], a local law was first derived for Wigner matrices. Subsequently in [7], a local law for the resolvent of  $(X - z)^\dagger(X - z)$  was established to prove the local circular law.

In generalizing the proof in [7] to our setting, a main difficulty is that the entries of  $TX$  are not independent. We will use a new comparison method proposed in [24], which roughly states that if the local laws hold for  $R(w)$  with Gaussian  $X$ , then they also hold in the case of a general  $X$ . For definiteness, we assume  $N = M$  for now, and  $T$  is a square matrix with singular decomposition  $T = UDV$ . For a Gaussian  $X \equiv X^{\text{Gauss}}$ , we have  $VX^{\text{Gauss}}U \stackrel{d}{=} \tilde{X}^{\text{Gauss}}$ , where  $\tilde{X}$  is another Gaussian random matrix. Then for the determinant in (1.1),

$$\det(TX^{\text{Gauss}} - z) = \det(DVX^{\text{Gauss}}U - z) \stackrel{d}{=} \det(D\tilde{X}^{\text{Gauss}} - z). \quad (1.3)$$

The problem is now reduced to the study of the singular values of  $D\tilde{X}^{\text{Gauss}} - z$ , which has independent entries. Notice the entries of  $D\tilde{X}^{\text{Gauss}}$  are not identically distributed, which will make our proof much more complicated. However, this issue can be handled, e.g. as in [14], where a local law was obtained for generalized Wigner matrices with non-identically distributed entries.

To use the comparison method invented in [24], it turns out the averaged local law from (1.2) is not sufficient. We have to control not only the trace of  $R(w)$ , but also the matrix  $R(w)$  itself by showing that  $R(w)$  is close to some deterministic matrix  $\Pi(w)$ , provided that  $\eta \gg N^{-1}$ . This closeness can be established in the sense of individual matrix entries  $R_{ij}(w) \approx \Pi_{ij}(w)$  (see e.g. [7, 17]). We call such an estimate an *entrywise local law*. More generally, in [4, 25] the following closeness was established for *generalized matrix entries*:

$$\langle \mathbf{v}, R(w)\mathbf{u} \rangle \approx \langle \mathbf{v}, \Pi(w)\mathbf{u} \rangle, \quad \eta \gg N^{-1}, \quad \forall \|\mathbf{v}\|_2, \|\mathbf{u}\|_2 = 1. \quad (1.4)$$

We call the estimate in (1.4) an *anisotropic local law*. (If  $\Pi$  is a scalar matrix, (1.4) is also referred to as an *isotropic local law*, in the sense that  $R(w)$  is approximately isotropic for large  $N$ .) This kind of anisotropic local law is needed in applying the method in [24]. Here we outline the three steps to establish the anisotropic local law for  $Q = (TX - z)^\dagger(TX - z)$ : (A) the entrywise local law and averaged local law when  $T$  is diagonal (Theorem 2.18); (B) the anisotropic local law when  $T$  is diagonal (Theorem 2.18); (C) the anisotropic local law and averaged local law when  $T$  is a general (rectangular) matrix (Theorem 2.19).

In performing Step (A), our proof is basically based on the methods in [7]. However, our multi-variable self-consistent equations and their solutions are much more complicated here. Thus a key part of the proof is to establish some basic properties of the asymptotic eigenvalue density and prove the stability of the self-consistent equations under small perturbations. These work need some new

ideas and analytic techniques (see Appendix A). In performing Step (B), we applied and extended the polynomialization method developed in [4, section 5]. Finally, as remarked around (1.3), (B) implies the anisotropic local law for a Gaussian  $X$  and a general  $T$ . Based on this fact we perform Step (C) using a self-consistent comparison argument in [24]. With the averaged local law proved in Step (C), we can prove the local circular law for  $TX$ . In general, the averaged local law we get is up to the non-optimal scale  $\eta \gg (N \wedge M)^{-1/2}$ . As a result, we can only prove the local circular law for  $TX$  up to the scale  $(N \wedge M)^{-1/4+\epsilon}$ . A new observation is that the non-optimal averaged local law can lead to the optimal local circular law for  $TX$  outside the unit circle (i.e.  $|z| > 1$ ) (see Section 2.4). To prove the optimal local circular law inside the unit circle (i.e.  $|z| < 1$ ), we need the optimal averaged local law up to the scale  $\eta \gg (N \wedge M)^{-1}$ , which can be obtained under the extra assumption that the entries of  $X$  have vanishing third moments.

**Conventions.** The fundamental large parameter is  $N$  and we assume that  $M$  is comparable to  $N$  (see (2.1)). All quantities that are not explicitly constant may depend on  $N$ , and we usually omit  $N$  from our notation. We use  $C$  to denote a generic large positive constant, which may depend on fixed parameters and whose value may change from one line to the next. Similarly, we use  $c$  or  $\epsilon$  to denote a generic small positive constant. If a constant depend on a quantity  $a$ , we use  $C(a)$  or  $C_a$  to indicate this dependence. We use  $\tau > 0$  in various assumptions to denote a small positive constant, and use  $\zeta, \tau'$  to denote constants that depend on  $\tau$  and may be chosen arbitrarily small. All constants  $C$ ,  $c$  and  $\epsilon$  may depend on  $\tau$ ; we neither indicate nor track this dependence.

For any (complex) matrix  $A$ , we use  $A^\dagger$  to denote its conjugate transpose,  $A^T$  the transpose,  $\|A\|$  the operator norm and  $\|A\|_{HS}$  the Hilbert-Schmidt norm. We use the notation  $\mathbf{v} = (v_i)_{i=1}^n$  for a vector in  $\mathbb{C}^n$ , and denote its Euclidean norm by  $|\mathbf{v}| \equiv \|\mathbf{v}\|_2$ . We usually write the  $n \times n$  identity matrix  $I_n$  as 1 without causing any confusions.

For two quantities  $A_N$  and  $B_N > 0$  depending on  $N$ , we use the notations  $A_N = O(B_N)$  and  $A_N \sim B_N$  to mean  $|A_N| \leq CB_N$  and  $C^{-1}B_N \leq |A_N| \leq CB_N$ , respectively, for some positive constant  $C > 0$ . We use  $A_N = o(B_N)$  to mean  $|A_N| \leq c_N B_N$  for some positive constant  $c_N \rightarrow 0$  as  $N \rightarrow \infty$ . If  $A_N$  is a matrix, we use the notations  $A_N = O(B_N)$  and  $A_N = o(B_N)$  to mean  $\|A_N\| = O(B_N)$  and  $\|A_N\| = o(B_N)$ , respectively.

**Acknowledgements.** The third author would like to thank Terence Tao, Mark Rudelson and Roman Vershynin for fruitful discussions and valuable suggestions.

## 2 The main results

In this section, we state and prove the main result of this paper. In Section 2.1, we define our model and list our main assumptions. In Section 2.2, we first define the asymptotic eigenvalue density  $\rho_{2c}$  of  $Q = (TX - z)^\dagger(TX - z)$ , and then state the main theorem—Theorem 2.6—of this paper. Its proof depends crucially on local estimates of the resolvent of  $Q$ , which are presented in Section 2.3. In Section 2.4, we prove Theorems 2.6 based on the local estimates stated in Section 2.3.

### 2.1 Definition of the model

In this paper, we want to understand the local statistics of the eigenvalues of  $TX - zI$ , where  $T$  is a deterministic  $N \times M$  matrix,  $X$  is a random  $M \times N$  matrix,  $z \in \mathbb{C}$  and  $I$  is the identity operator. We assume  $M \sim N$ , i.e.

$$\tau \leq \frac{M}{N} \leq \tau^{-1} \quad (2.1)$$

for some small  $\tau > 0$ . We assume the entries  $X_{i\mu}$  of  $X$  are independent (not necessarily identically distributed) random variables satisfying

$$\mathbb{E} X_{i\mu} = 0, \quad \mathbb{E} |X_{i\mu}|^2 = \frac{1}{N \wedge M} \quad (2.2)$$

for all  $1 \leq i \leq M, 1 \leq \mu \leq N$ . For definiteness, in this paper we only focus on the case where all matrix entries are real. However, our results and proofs also hold, after minor changes, in the complex case if we assume in addition  $\mathbb{E} X_{i\mu}^2 = 0$  for  $X_{i\mu} \in \mathbb{C}$ . We assume that for all  $p \in \mathbb{N}$ , there is an  $N$ -independent constant  $C_p$  such that

$$\mathbb{E} |\sqrt{N \wedge M} X_{i\mu}|^p \leq C_p \quad (2.3)$$

for all  $1 \leq i \leq M, 1 \leq \mu \leq N$ . We define  $\Sigma := TT^\dagger$ , and assume the eigenvalues of  $\Sigma$  satisfy that

$$\tau^{-1} \geq \sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{N \wedge M} \geq \tau \quad (2.4)$$

and all other eigenvalues are 0. We can normalize  $T$  by multiplying a scalar such that

$$\frac{1}{N \wedge M} \sum_{i=1}^{N \wedge M} \sigma_i = 1. \quad (2.5)$$

We summarize our basic assumptions here for future reference.

**Assumption 2.1.** *We suppose that (2.1), (2.2), (2.3), (2.4) and (2.5) hold.*

## 2.2 The main theorem

Our main result is Theorem 2.6. To state it, we need to define the asymptotic eigenvalue density function for  $Q$ . We first introduce the self-consistent equations, and the asymptotic eigenvalue density will be closely related to their solutions. Define

$$\rho_\Sigma := \frac{1}{N \wedge M} \sum_{i=1}^{N \wedge M} \delta_{\sigma_i} \quad (2.6)$$

as the empirical spectral density of  $\Sigma$ . Let  $n := |\text{supp } \rho_\Sigma|$  be the number of distinct nonzero eigenvalues of  $\Sigma$ , which are denoted as

$$\tau^{-1} \geq s_1 > s_2 > \cdots > s_n \geq \tau. \quad (2.7)$$

Let  $l_i$  be the multiplicity of  $s_i$ . By (2.5),  $l_i$  and  $s_i$  satisfy the normalization conditions

$$\frac{1}{N \wedge M} \sum_{i=1}^n l_i = 1, \quad \frac{1}{N \wedge M} \sum_{i=1}^n l_i s_i = 1. \quad (2.8)$$

For each  $w \in \mathbb{C}_+ := \{w \in \mathbb{C} : \text{Im } w > 0\}$ , we define the self-consistent equations of  $(m_1, m_2)$  as

$$\frac{1}{m_2} = -w(1 + m_1) + \frac{|z|^2}{1 + m_1}, \quad (2.9)$$

$$m_1 = \frac{1}{N} \sum_{i=1}^n l_i s_i \left[ -w(1 + s_i m_2) + \frac{|z|^2}{1 + m_1} \right]^{-1}. \quad (2.10)$$

If we plug (2.9) into (2.10), we get the self-consistent equation for  $m_1$  only,

$$m_1 = \frac{1}{N} \sum_{i=1}^n l_i s_i \left[ -w \left( 1 + \frac{s_i}{-w(1+m_1) + \frac{|z|^2}{1+m_1}} \right) + \frac{|z|^2}{1+m_1} \right]^{-1}. \quad (2.11)$$

The next lemma states that the solution to (2.11) in  $\mathbb{C}_+$  is unique if  $z$  is away from the unit circle. It is proved in Appendix A.3.

**Lemma 2.2.** *Fix  $z \in \mathbb{C}$  such that  $|z| \neq 1$ . For  $w \in \mathbb{C}_+$ , there exists at most one analytic function  $m_{1c,z,\Sigma}(w) : \mathbb{C}_+ \rightarrow \mathbb{C}_+$  such that (2.11) holds and  $wm_{1c,z,\Sigma}(w) \in \mathbb{C}_+$ . Moreover,  $m_{1c,z,\Sigma,N}(w)$  is the Stieltjes transform of a positive integrable function  $\rho_{1c}$  with compact support in  $[0, \infty)$ .*

We shall abbreviate  $m_{1c}(w) := m_{1c,z,\Sigma}(w)$ . We also define  $m_{2c}(w) := m_{2c,z,\Sigma}(w)$  by taking  $m_1 = m_{1c}(w)$  in (2.9). Obviously,  $m_{2c}$  is also an analytic function of  $w$ . Furthermore, for any  $w \in \mathbb{C}_+$  we have  $m_{2c}(w), wm_{2c}(w) \in \mathbb{C}_+$  by using (2.9) and  $m_{1c}, wm_{1c} \in \mathbb{C}_+$ . We define two functions on  $\mathbb{R}$  as

$$\rho_{1,2c}(x) = \frac{1}{\pi} \lim_{\eta \searrow 0} \text{Im } m_{1,2c}(x + i\eta), \quad x \in \mathbb{R}. \quad (2.12)$$

It is easy to see that  $\rho_{1,2c} \geq 0$  and  $\text{supp}(\rho_{1,2c}) \subseteq [0, \infty)$ . Moreover,  $\text{supp } \rho_{2c} = \text{supp } \rho_{1c}$  by (2.9). We shall call  $\rho_{2c}$  the asymptotic eigenvalue density of  $Q = (TX - z)^\dagger (TX - z)$  (for a reason that will be made clear during the proof in Section 4). Since  $\text{Im}(wm_{2c}) \geq 0$ , we have

$$\text{Im} \left[ -w(1 + s_i m_{2c}) + \frac{|z|^2}{1 + m_{1c}} \right] \leq -\text{Im } w,$$

and (2.10) gives  $|m_{1c}| \leq 1/\text{Im } w \rightarrow 0$  as  $\text{Im } w \rightarrow \infty$ . Similarly,  $|m_{2c}| \leq 1/\text{Im } w \rightarrow 0$  as  $\text{Im } w \rightarrow \infty$ . Thus  $m_{1,2c}(w)$  is indeed the Stieltjes transform of  $\rho_{1,2c}$ ,

$$m_{1,2c}(w) = \int_{\mathbb{R}} \frac{\rho_{1,2c}(x)}{x - w} dx. \quad (2.13)$$

We now state the basic properties of  $\rho_{1c}$  and  $\rho_{2c}$ , which can be obtained by studying the solutions  $m_{1,2c}(w)$  to the self-consistent equations (2.9) and (2.11) when  $w \in (0, \infty)$ . Here we extend the definition of  $m_{1,2c}$  continuously down to the real axis by setting

$$m_{1,2c}(x) = \lim_{\eta \searrow 0} m_{1,2c}(x + i\eta), \quad x \in \mathbb{R}.$$

As a convention, for  $w \in \overline{\mathbb{C}_+}$ , we take  $\sqrt{w}$  to be the branch with positive imaginary part. Define  $m := \sqrt{w}(1 + m_1)$  and  $m_c := \sqrt{w}(1 + m_{1c})$ . Equation (2.11) then becomes

$$f(\sqrt{w}, m) = 0, \quad (2.14)$$

where

$$f(\sqrt{w}, m) = -\sqrt{w} + m + \frac{1}{N} \sum_{i=1}^n l_i s_i \frac{m(m^2 - |z|^2)}{\sqrt{w}m^3 - (s_i + |z|^2)m^2 - \sqrt{w}|z|^2m + |z|^4}. \quad (2.15)$$

The following lemma gives the basic structure of  $\text{supp } \rho_{1,2c}$ . Its proof is given in Appendix A.1.

**Lemma 2.3.** Fix  $\tau \leq ||z|^2 - 1| \leq \tau^{-1}$ . The support of  $\rho_{1,2c}$  is a union of connected components:

$$\text{supp } \rho_{1,2c} \cap (0, +\infty) = \left( \bigcup_{1 \leq k \leq L} [e_{2k}, e_{2k-1}] \right) \cap (0, \infty), \quad (2.16)$$

where  $L \equiv L(n) \in \mathbb{N}$  and  $C_1 \tau^{-1} \geq e_1 > e_2 > \dots > e_{2L} \geq 0$  for some constant  $C_1 > 0$  that does not depend on  $\tau$ . If  $|z|^2 \leq 1 - \tau$ , we have  $e_{2L} = 0$ ; if  $1 + \tau \leq |z|^2 \leq 1 + \tau^{-1}$ ,  $e_{2L} \geq \epsilon(\tau)$  for some constant  $\epsilon(\tau) > 0$ . Moreover, for every  $e_i > 0$ , there exists a unique  $m_c(e_i)$  such that

$$\partial_m f(\sqrt{e_i}, m_c(e_i)) = 0. \quad (2.17)$$

We shall call  $e_i$ 's the edges of  $\rho_{1c}$ . For any  $w \in (0, \infty)$  and  $1 \leq i \leq n$ , the cubic polynomial  $\sqrt{w}m^3 - (s_i + |z|^2)m^2 - \sqrt{w}|z|^2m + |z|^4$  in (2.15) has three distinct roots  $a_i(w) > 0$ ,  $b_i(w) > 0$  and  $-c_i(w) < 0$  (see Lemma A.1). Our next assumption on  $\rho_\Sigma$  and  $|z|$  takes the form of the following regularity conditions.

**Definition 2.4.** (Regularity) Fix  $\tau \leq ||z|^2 - 1| \leq \tau^{-1}$  and a small constant  $\epsilon > 0$ .

(i) We say that the edge  $e_k \neq 0$ ,  $k = 1, \dots, 2L$ , is regular if

$$\min_{1 \leq i \leq n} \{|m_c(e_k) - a_i(e_k)|, |m_c(e_k) - b_i(e_k)|, |m_c(e_k) + c_i(e_k)|\} \geq \epsilon, \quad (2.18)$$

and

$$|\partial_m^2 f(\sqrt{e_k}, m_c(e_k))| \geq \epsilon. \quad (2.19)$$

In the case  $|z|^2 \leq 1 - \tau$ , we always call  $e_{2L} = 0$  a regular edge.

(ii) We say that the bulk components  $[e_{2k}, e_{2k-1}]$  is regular if for any fixed  $\tau' > 0$  there exists a constant  $c(\tau, \tau') > 0$  such that the density of  $\rho_{1c}$  in  $[e_{2k} + \tau', e_{2k-1} - \tau']$  is bounded from below by  $c$ .

*Remark 1:* The edge regularity conditions (i) has previously appeared (may be in slightly different forms) in several works on sample covariance matrices and Wigner matrices [3, 11, 23, 24, 26, 29]. The conditions (2.18) and (2.19) guarantees a regular square-root behavior of  $\rho_{1c}$  near  $e_k$  and ensures that the gap in the spectrum of  $\rho_{1c}$  adjacent to  $e_k$  does not close for large  $N$  (Lemma A.5),

$$\min_{l \neq k} |e_l - e_k| \geq \epsilon \quad (2.20)$$

for some constant  $\epsilon > 0$ . The bulk regularity condition (ii) was introduced in [24]. It imposes a lower bound on the density of eigenvalues away from the edges. Without it, one can have points in the interior of  $\text{supp } \rho_{1c}$  with an arbitrarily small density and our arguments would fail.

*Remark 2:* The regularity conditions in Definition 2.4 are stable under perturbations of  $|z|$  and  $\rho_\Sigma$ . In particular, fix  $\rho_\Sigma$ , suppose the regularity conditions are satisfied at  $z = z_0$  with  $\tau \leq ||z_0|^2 - 1| \leq \tau^{-1}$ . Then for sufficiently small  $c > 0$ , the regularity conditions hold uniformly in  $z \in \{z : ||z| - |z_0|| \leq c\}$ . For a detailed discussion, see the remark at the end of Section A.3.

We will use the following notion of stochastic domination, which was first introduced in [12] and subsequently used in many works on random matrix theory, such as [4, 5, 7, 13, 14, 24]. It simplifies the presentation of the results and their proofs by systematizing statements of the form “ $\xi$  is bounded by  $\zeta$  with high probability up to a small power of  $N$ ”.

**Definition 2.5** (Stochastic domination). (i) Let

$$\xi = \left( \xi^{(N)}(u) : N \in \mathbb{N}, u \in U^{(N)} \right), \quad \zeta = \left( \zeta^{(N)}(u) : N \in \mathbb{N}, u \in U^{(N)} \right)$$

be two families of nonnegative random variables, where  $U^{(N)}$  is a possibly  $N$ -dependent parameter set. We say  $\xi$  is stochastically dominated by  $\zeta$ , uniformly in  $u$ , if for any (small)  $\epsilon > 0$  and (large)  $D > 0$ ,

$$\sup_{u \in U^{(N)}} \mathbb{P} \left[ \xi^{(N)}(u) > N^\epsilon \zeta^{(N)}(u) \right] \leq N^{-D}$$

for large enough  $N \geq N_0(\epsilon, D)$ , and we use the notation  $\xi < \zeta$ . Throughout this paper the stochastic domination will always be uniform in all parameters that are not explicitly fixed (such as matrix indices, and  $w$  and  $z$  that take values in some compact sets). Note that  $N_0(\epsilon, D)$  may depend on quantities that are explicitly constant, such as  $\tau$  and  $C_p$  in (2.1), (2.3) and (2.4).

(ii) If for some complex family  $\xi$  we have  $|\xi| < \zeta$ , we also write  $\xi < \zeta$  or  $\xi = O_{<}(\zeta)$ . We also extend the definition of  $O_{<}(\cdot)$  to matrices in the weak operator sense as follows. Let  $A$  be a family of complex square random matrices and  $\zeta$  a family of nonnegative random variables. Then we use  $A = O_{<}(\zeta)$  to mean  $\|A\| < \zeta$ , where  $\|A\|$  is the operator norm of  $A$ .

(iv) We say that an event  $\Xi$  holds with high probability if  $1 - 1(\Xi) < 0$ .

In the following, we denote the eigenvalues of  $TX$  as  $\mu_j$ ,  $1 \leq j \leq N$ . We are now ready to state our main theorem, i.e. the general local circular law for  $TX$ .

**Theorem 2.6** (Local circular law for  $TX$ ). *Suppose Assumption 2.1 holds, and  $\tau \leq ||z_0|^2 - 1| \leq \tau^{-1}$  for any  $N$  ( $z_0$  can depend on  $N$ ). Suppose  $\rho_\Sigma$  (defined in (2.6)) and  $|z_0|$  are such that all the edges and bulk components of  $\rho_{1c}$  are regular in the sense of Definition 2.4. We assume in addition that the entries of  $X$  have a density bounded by  $N^{C_2}$  for some  $C_2 > 0$ . Let  $F$  be a smooth non-negative function which may depend on  $N$ , such that  $\|F\|_\infty \leq C_1$ ,  $\|F'\|_\infty \leq N^{C_1}$  and  $F(z) = 0$  for  $|z| \geq C_1$ , for some constant  $C_1 > 0$  independent of  $N$ . Let  $F_{z_0,a}(z) = K^{2a} F(K^a(z - z_0))$ , where  $K := N \wedge M$ . Then  $TX$  has  $(N - K)$  trivial zero eigenvalues, and for the other eigenvalues  $\mu_j$ ,  $1 \leq j \leq K$ , we have*

$$\frac{1}{K} \sum_{j=1}^K F_{z_0,a}(\mu_j) - \frac{1}{\pi} \int F_{z_0,a}(z) \tilde{\chi}_{\mathbb{D}}(z) dA(z) < K^{-1/2+2a} \|\Delta F\|_{L^1}, \quad (2.21)$$

for any  $a \in (0, 1/4]$ . Here

$$\tilde{\chi}_{\mathbb{D}}(z) := \frac{1}{4} \int_0^\infty (\log x) \Delta_z \rho_{2c}(x, z) dx, \quad (2.22)$$

where  $\rho_{2c} \equiv \rho_{2c,z,\Sigma}$  is defined in (2.12). If  $1 + \tau \leq |z_0|^2 \leq 1 + \tau^{-1}$  or the entries of  $X$  have vanishing third moments,

$$\mathbb{E} X_{i\mu}^3 = 0, \quad (2.23)$$

for  $1 \leq i \leq M, 1 \leq \mu \leq N$ , then we have the improved result

$$\frac{1}{K} \sum_{j=1}^K F_{z_0,a}(\mu_j) - \frac{1}{\pi} \int F_{z_0,a}(z) \tilde{\chi}_{\mathbb{D}}(z) dA(z) < K^{-1+2a} \|\Delta F\|_{L^1}, \quad (2.24)$$

for any  $a \in (0, 1/2]$ . If  $N = M$ , the bounded density condition for the entries of  $X$  is not necessary.

*Remark 1:* Note that  $F_{z_0,a}(z) = K^{2a} F(K^a(z - z_0))$  is an approximate delta function obtained from rescaling  $F$  to the size of order  $K^{-a}$  around  $z_0$ . Thus (2.21) gives the general circular law up to scale  $K^{-1/4+\epsilon}$ , while (2.24) gives the general circular law up to scale  $K^{-1/2+\epsilon}$ . The  $\tilde{\chi}_{\mathbb{D}}$  in (2.22) gives the distribution of the eigenvalues of  $TX$ . It is rotationally symmetric, because  $\rho_{2c}(x, z)$  only depends on  $|z|$  (see (2.9) and (2.10)). When  $T$  is the identity matrix,  $\tilde{\chi}_{\mathbb{D}}$  becomes the indicator function  $\chi_{\mathbb{D}}$  on the unit disk  $\mathbb{D}$ , and we get the well-known local circular law for  $X$  [7]. For a general  $T$ , we do



not have much understanding of  $\tilde{\chi}_{\mathbb{D}}$  so far. This will be one of the topics of our future study. Also, we have assumed that  $z$  is strictly away from the unit circle. Our proof may be extended to the  $|z - 1| = o(1)$  case if we have a better understanding of the solutions  $m_{1,2c}$  to equations (2.9) and (2.10).

*Remark 2:* As explained in the Introduction, the basic strategy of this paper is first to prove the anisotropic local law for the resolvent of  $Q$  when  $X$  is Gaussian, and then to get the anisotropic local law for a general  $X$  through comparison with the Gaussian case. Without (2.23), our comparison arguments do not give the anisotropic local law up to the optimal scale, so we can only prove the weaker bound (2.21). We will try to remove this assumption in future works.

*Remark 3:* In the statement of the theorem, we have included an extra bounded density condition. This is only used in Lemma 2.23 to give a lower bound for the smallest singular value of  $TX - z$ . Thus it can be removed if we have a stronger result about the smallest singular value.

We conclude this section with two examples verifying the regularity conditions of Definition 2.4.

**Example 2.7** (Bounded number of distinct eigenvalues). *We suppose that  $n$  is fixed, and that  $s_1, \dots, s_n$  and  $\rho_{\Sigma}(\{s_1\}), \dots, \rho_{\Sigma}(\{s_n\})$  all converge as  $N \rightarrow \infty$ . We suppose that  $\lim_N e_k > \lim_N e_{k+1}$  for all  $k$ , and furthermore for all  $e_k$  we have  $\partial_m^2 f(\sqrt{e_k}, m_c(e_k)) \neq 0$ . Then it is easy to check that all the edges and bulk components are regular in the sense of Definition 2.4 for small enough  $\epsilon$ .*

**Example 2.8** (Continuous limit). *We suppose  $\rho_{\Sigma}$  is supported in some interval  $[a, b] \subset (0, \infty)$ , and that  $\rho_{\Sigma}$  converges in distribution to some measure  $\rho_{\infty}$  that is absolutely continuous and whose density satisfies  $\tau \leq d\rho_{\infty}(E)/dE \leq \tau^{-1}$  for  $E \in [a, b]$ . Then there are only a small number (which is independent of  $n$ ) of connected components for  $\text{supp } \rho_{1c}$ , and all the edges and bulk components are regular. See the remark at the end of Section A.1.*

## 2.3 Hermitization and local laws for resolvents

In the following, we use the notation

$$Y \equiv Y_z := TX - zI, \quad (2.25)$$

where  $I$  is the identity matrix. Following Girko's Hermitization technique [20], the first step in proving the local circular law is to understand the local statistics of singular values of  $Y$ . In this subsection, we present the main local estimates concerning the resolvents  $(YY^\dagger - w)^{-1}$  and  $(Y^\dagger Y - w)^{-1}$ . These results will be used later to prove Theorem 2.6.

Our local laws can be formulated in a simple, unified fashion using a  $2N \times 2N$  block matrix, which is a linear function of  $X$ .

**Definition 2.9** (Index sets). *We define the index sets*

$$\mathcal{I}_1 := \{1, \dots, N\}, \quad \mathcal{I}_1^M := \{1, \dots, M\}, \quad \mathcal{I}_2 := \{N+1, \dots, 2N\}, \quad \mathcal{I} := \mathcal{I}_1 \cup \mathcal{I}_2, \quad \mathcal{I}^M := \mathcal{I}_1^M \cup \mathcal{I}_2.$$

*We will consistently use the latin letters  $i, j \in \mathcal{I}_1$  or  $\mathcal{I}_1^M$ , greek letters  $\mu, \nu \in \mathcal{I}_2$ , and  $s, t \in \mathcal{I}$ . We label the indices of the matrices according to*

$$X = (X_{i\mu} : i \in \mathcal{I}_1^M, \mu \in \mathcal{I}_2), \quad T = (T_{ij} : i \in \mathcal{I}_1, j \in \mathcal{I}_1^M).$$

*When  $M = N$ , we always identify  $\mathcal{I}_1^M$  with  $\mathcal{I}_1$ . For  $i \in \mathcal{I}_1$  and  $\mu \in \mathcal{I}_2$ , we introduce the notations  $\bar{i} := i + N \in \mathcal{I}_2$  and  $\bar{\mu} := \mu - N \in \mathcal{I}_1$ .*

**Definition 2.10** (Groups). For an  $\mathcal{I} \times \mathcal{I}$  matrix  $A$ , we define the  $2 \times 2$  matrix  $A_{[ij]}$  as

$$A_{[ij]} = \begin{pmatrix} A_{ij} & A_{i\bar{j}} \\ A_{\bar{i}j} & A_{\bar{i}\bar{j}} \end{pmatrix}. \quad (2.26)$$

We shall call  $A_{[ij]}$  a diagonal group if  $i = j$ , and an off-diagonal group otherwise.

**Definition 2.11** (Linearizing block matrix). For  $w := E + i\eta \in \mathbb{C}_+$ , we define the  $\mathcal{I} \times \mathcal{I}$  matrix

$$H(w) \equiv H(T, X, z, w) := \begin{pmatrix} -wI & w^{1/2}Y \\ w^{1/2}Y^\dagger & -wI \end{pmatrix}, \quad (2.27)$$

where we take the branch of  $\sqrt{w}$  with positive imaginary part. Define the  $\mathcal{I} \times \mathcal{I}$  matrix

$$G(w) \equiv G(T, X, z, w) := H(w)^{-1}, \quad (2.28)$$

as well as the  $\mathcal{I}_1 \times \mathcal{I}_1$  and  $\mathcal{I}_2 \times \mathcal{I}_2$  matrices

$$G_L(w) = (Y Y^\dagger - w)^{-1}, \quad G_R(w) = (Y^\dagger Y - w)^{-1}. \quad (2.29)$$

Throughout the following, we frequently omit the argument  $w$  from our notations.

By Schur's complement formula, it is easy to see that

$$G(w) = \begin{pmatrix} G_L & w^{-1/2}G_L Y \\ w^{-1/2}Y^\dagger G_L & w^{-1}Y^\dagger G_L Y - w^{-1}I \end{pmatrix} = \begin{pmatrix} w^{-1}Y G_R Y^\dagger - w^{-1}I & w^{-1/2}Y G_R \\ w^{-1/2}G_R Y^\dagger & G_R \end{pmatrix}. \quad (2.30)$$

Therefore a control of  $G$  immediately yields controls of the resolvents  $G_L$  and  $G_R$ .

In the following, we only consider the  $N \leq M$  case. The  $N > M$  case, as we will see, will be built easily upon  $N \leq M$  case. We introduce a deterministic matrix  $\Pi$ , which will be proved to be close to  $G$  with high probability.

**Definition 2.12** (Deterministic limit of  $G$ ). Suppose  $N \leq M$  and  $T$  has a singular decomposition

$$T = U \bar{D} V, \quad \bar{D} = (D, 0), \quad (2.31)$$

where  $D = \text{diag}(d_1, d_2, \dots, d_N)$  is a diagonal matrix. Define  $\pi_{[i]c}$  to be the  $2 \times 2$  matrix such that

$$(\pi_{[i]c})^{-1} = \begin{pmatrix} -w(1 + |d_i|^2 m_{2c}) & -w^{1/2}z \\ -w^{1/2}\bar{z} & -w(1 + m_{1c}) \end{pmatrix}. \quad (2.32)$$

Let  $\Pi_d$  be the  $2N \times 2N$  matrix with  $(\Pi_d)_{[ii]} = \pi_{[i]c}$  and all other entries being zero. Define

$$\Pi \equiv \Pi(\Sigma, z, w) := \begin{pmatrix} U & 0 \\ 0 & U \end{pmatrix} \Pi_d \begin{pmatrix} U^\dagger & 0 \\ 0 & U^\dagger \end{pmatrix} = \begin{pmatrix} -(1 + m_{1c})A(\Sigma) & w^{-1/2}zA(\Sigma) \\ w^{-1/2}\bar{z}A(\Sigma) & -(1 + m_{2c}\Sigma)A(\Sigma) \end{pmatrix}, \quad (2.33)$$

where  $\Sigma = TT^\dagger$  and  $A(\Sigma) = [w(1 + m_{2c}\Sigma)(1 + m_{1c}) - |z|^2]^{-1}$ .

**Definition 2.13** (Averaged variables). Suppose  $N \leq M$ . Define the averaged random variables

$$m_1 := \frac{1}{N} \sum_{i \in \mathcal{I}_1} (\bar{\Sigma} G)_{ii}, \quad m_2 := \frac{1}{N} \sum_{\mu \in \mathcal{I}_2} (\bar{\Sigma} G)_{\mu\mu}, \quad (2.34)$$

where

$$\bar{\Sigma} := \begin{pmatrix} \Sigma & 0 \\ 0 & I \end{pmatrix}. \quad (2.35)$$

Define  $\pi_{[i]}$  to be the  $2 \times 2$  matrix such that

$$(\pi_{[i]})^{-1} = \begin{pmatrix} -w(1 + |d_i|^2 m_2) & -w^{1/2} z \\ -w^{1/2} \bar{z} & -w(1 + m_1) \end{pmatrix}. \quad (2.36)$$

*Remark:* Note that under the above definition we have

$$m_2 = \frac{1}{N} \text{Tr } G_R = \frac{1}{N} \text{Tr } G_L,$$

which is the Stieltjes transform of the empirical eigenvalue density of  $YY^\dagger$  and  $Y^\dagger Y$ . Moreover, we will see from the proof that  $m_{1,2c}$  are the almost sure limits of  $m_{1,2}$  as  $N \rightarrow \infty$  with

$$m_{1c} = \frac{1}{N} \sum_{i \in \mathcal{I}_1} (\bar{\Sigma} \Pi)_{ii}, \quad m_{2c} = \frac{1}{N} \sum_{\mu \in \mathcal{I}_2} (\bar{\Sigma} \Pi)_{\mu\mu}. \quad (2.37)$$

The following two propositions summarize the properties of  $\rho_{1,2c}$  and  $m_{1,2c}$  that are needed to understand the main results in this section. They are proved in Appendix A. In Fig. 2 we plot  $\rho_{2c}$  for the example from Fig. 1 in the cases  $|z| > 1$  and  $|z| < 1$ , respectively.

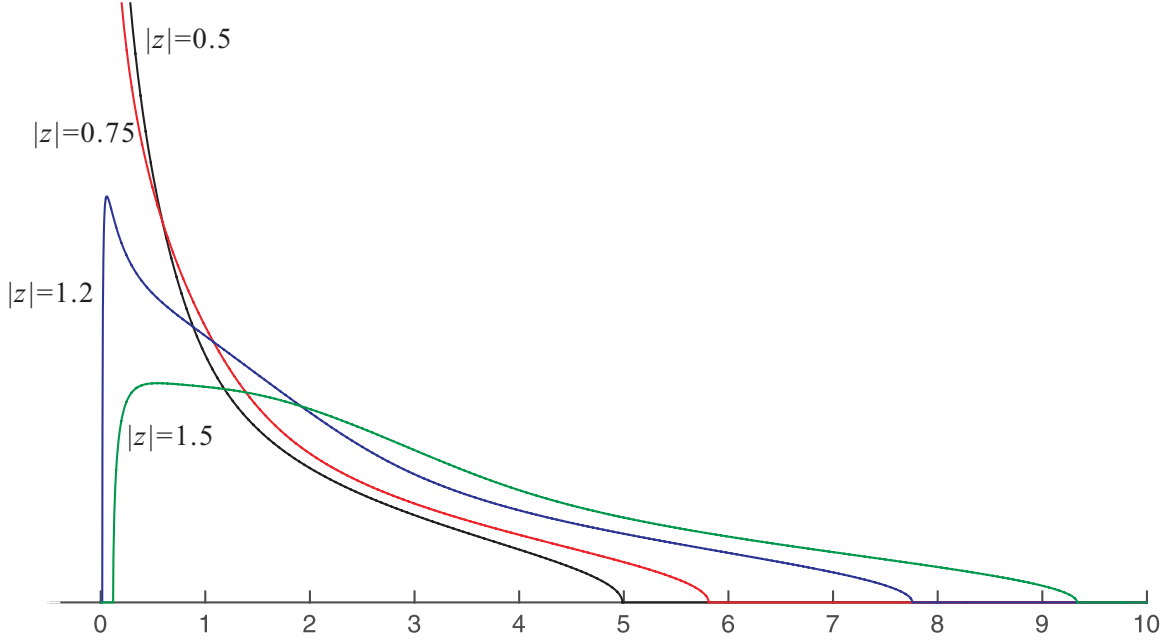


Figure 2: The densities  $\rho_{2c}(x, z)$  when  $|z| = 0.5, 0.75, 1.2, 1.5$ . Here  $\rho_\Sigma = 0.5\delta_{\sqrt{2/17}} + 0.5\delta_{4\sqrt{2/17}}$ .

**Proposition 2.14** (Basic properties of  $\rho_{1,2c}$ ). *Fix  $\epsilon > 0$ . The density  $\rho_{1c}$  is compactly supported in  $[0, \infty)$  and the following properties regarding  $\rho_{1c}$  hold.*

(i) The support of  $\rho_{1c}$  is  $\bigcup_{1 \leq k \leq L(n)} [e_{2k}, e_{2k-1}]$  where  $e_1 > e_2 > \dots > e_{2L} \geq 0$ . If  $1 + \tau \leq |z|^2 \leq 1 + \tau^{-1}$ , then  $e_{2L} \geq \epsilon$ ; if  $|z|^2 \leq 1 - \tau$ , then  $e_1 = 0$ .

(ii) Suppose  $[e_{2k}, e_{2k-1}]$  is a regular bulk component. For any  $\tau' > 0$ , if  $x \in [e_{2k} + \tau', e_{2k-1} - \tau']$ , then  $\rho_{1c}(x) \sim 1$ .

(iii) Suppose  $e_j$  is a nonzero regular edge. If  $j$  is even, then  $\rho_{1c}(x) \sim \sqrt{x - e_j}$  as  $x \rightarrow e_j$  from above. Otherwise if  $j$  is odd, then  $\rho_{1c}(x) \sim \sqrt{e_j - x}$  as  $x \rightarrow e_j$  from below.

(iv) If  $|z|^2 \leq 1 - \tau$ , then  $\rho_{1c}(x) \sim x^{-1/2}$  as  $x \searrow e_{2L} = 0$ .

The same results also hold for  $\rho_{2c}$ . In addition,  $\rho_{2c}$  is a probability density.

**Proposition 2.15.** *The preceding proposition implies that, uniformly in  $w$  in any compact set of  $\mathbb{C}_+$ ,*

$$|m_{1,2c}(w)| = O(|w|^{-1/2}). \quad (2.38)$$

Moreover, if  $1 + \tau \leq |z|^2 \leq 1 + \tau^{-1}$ , then  $|m_{1,2c}(w)| \sim 1$  for  $w$  in any compact set of  $\mathbb{C}_+$ ; if  $|z|^2 \leq 1 - \tau$ , then  $|m_{1,2c}(w)| \sim |w|^{-1/2}$  for  $w$  in any compact set of  $\mathbb{C}_+$ .

We will consistently use the notation  $E + i\eta$  for the spectral parameter  $w$ . In this paper, we regard the quantities  $E(w)$  and  $\eta(w)$  as functions of  $w$  and usually omit the argument  $w$ . In the following we would like to define several spectral domains of  $w$  that will be used in the proof.

**Definition 2.16** (Spectral domains). *Fix a small constant  $\zeta > 0$  which may depend on  $\tau$ . The spectral parameter  $w$  is always assumed to be in the fundamental domain*

$$\mathbf{D} \equiv \mathbf{D}(\zeta, N) := \{w \in \mathbb{C}_+ : 0 \leq E \leq \zeta^{-1}, N^{-1+\zeta}|m_{2c}|^{-1} \leq \eta \leq \zeta^{-1}\}. \quad (2.39)$$

unless otherwise indicated. Given a regular edge  $e_k$ , we define the subdomain

$$\mathbf{D}_k^e \equiv \mathbf{D}_k^e(\zeta, \tau', N) := \{w \in \mathbf{D}(\zeta, N) : |E - e_k| \leq \tau', E \geq 0\}. \quad (2.40)$$

Corresponding to a regular bulk component  $[e_{2k}, e_{2k-1}]$ , we define the subdomain

$$\mathbf{D}_k^b \equiv \mathbf{D}_k^b(\zeta, \tau', N) := \{w \in \mathbf{D}(\zeta, N) : E \in [e_{2k} + \tau', e_{2k-1} - \tau']\}. \quad (2.41)$$

For the component outside  $\text{supp } \rho_{1c}$ , we define the subdomain

$$\mathbf{D}^o \equiv \mathbf{D}^o(\zeta, \tau', N) := \{w \in \mathbf{D}(\zeta, N) : \text{dist}(E, \text{supp } \rho_{1c}) \geq \tau'\}. \quad (2.42)$$

We also need the following domain with large  $\eta$ ,

$$\mathbf{D}_L \equiv \mathbf{D}_L(\zeta) := \{w \in \mathbb{C}_+ : 0 \leq E \leq \zeta^{-1}, \eta \geq \zeta^{-1}\}, \quad (2.43)$$

and the subdomain of  $\mathbf{D} \cup \mathbf{D}_L$ ,

$$\hat{\mathbf{D}} \equiv \hat{\mathbf{D}}(\zeta, N) := \{w \in \mathbf{D}(\zeta, N) : \eta \geq N^{-1/2+\zeta}|m_{2c}|^{-1}\} \cup \mathbf{D}_L(\zeta). \quad (2.44)$$

We call  $\mathbf{S}$  a regular domain if it is a regular  $\mathbf{D}_k^e$  or  $\mathbf{D}_k^b$  domain, a  $\mathbf{D}^o$  domain or a  $\mathbf{D}_L$  domain.

*Remark:* In the definition of  $\mathbf{D}$ , we have suppressed the explicit  $w$ -dependence. Notice that when  $|z|^2 < 1 - \tau$ , since  $|m_{2c}| \sim |w|^{-1/2}$  as  $w \rightarrow 0$ , we allow  $\eta \sim |w| \sim N^{-2+2\zeta}$  in  $\mathbf{D}$ . In the definition of  $\mathbf{D}_k^e$ , the condition  $E \geq 0$  is only for the edge at 0 when  $|z|^2 \leq 1 - \tau$ .

Now we are prepared to state the various local laws satisfied by  $G$  defined in (2.28). Let

$$\Psi \equiv \Psi(w) := \sqrt{\frac{\text{Im}(m_{1c} + m_{2c})}{N\eta}} + \frac{1}{N\eta} \quad (2.45)$$

be the deterministic control parameter.

**Definition 2.17** (Local laws). Suppose  $N \leq M$ . Recall  $G \equiv G(T, X, z, w)$  defined in (2.28) and  $\Pi \equiv \Pi(\Sigma, z, w)$  defined in (2.33). Let  $\mathbf{S}$  be a regular domain.

(i) We say that the entrywise local law holds with parameters  $(T, X, z, \mathbf{S})$  if

$$[G(T, X, z, w) - \Pi(\Sigma, z, w)]_{st} < \Psi(w) \quad (2.46)$$

uniformly in  $w \in \mathbf{S}$  and  $s, t \in \mathcal{I}$ .

(ii) We say that the anisotropic local law holds with parameters  $(T, X, z, \mathbf{S})$  if

$$\|G(T, X, z, w) - \Pi(\Sigma, z, w)\| < \Psi(w) \quad (2.47)$$

uniformly in  $w \in \mathbf{S}$ .

(iii) We say that the averaged local law holds with parameters  $(T, X, z, \mathbf{S})$  if

$$|m_2(T, X, z, w) - m_{2c}(\Sigma, z, w)| < \frac{1}{N\eta} \quad (2.48)$$

uniformly in  $w \in \mathbf{S}$ .

The local laws for  $G$  with a general  $T$  will be built upon the following result with a diagonal  $T$ .

**Theorem 2.18** (Local laws when  $T$  is diagonal). Fix  $\tau \leq ||z|^2 - 1| \leq \tau^{-1}$ . Suppose Assumption 2.1 holds,  $N = M$ , and  $T \equiv D := \text{diag}(d_1, \dots, d_N)$  is a diagonal matrix. Let  $\mathbf{S}$  be a regular domain. Then the entrywise local law, anisotropic local law and averaged local law hold with parameters  $(D, X, z, \mathbf{S})$ .

Now suppose that  $N \leq M$  and  $T$  is an  $N \times M$  matrix such that the eigenvalues of  $\Sigma$  satisfy (2.4) and (2.5). Consider the singular decomposition  $T = U\bar{D}V$ , where  $U$  is an  $N \times N$  unitary matrix,  $V$  is an  $M \times M$  unitary matrix and  $\bar{D} = (D, 0)$  is an  $N \times M$  matrix such that  $D = \text{diag}(d_1, d_2, \dots, d_N)$ . Then we have

$$TX - z = UDV_1X - z, \quad (2.49)$$

where  $V_1$  is an  $N \times M$  matrix and  $V_2$  is an  $(M - N) \times M$  matrix defined through  $V = \begin{pmatrix} V_1 \\ V_2 \end{pmatrix}$ . If

$X = X^{\text{Gauss}}$  is Gaussian, then  $V_1X^{\text{Gauss}} \stackrel{d}{=} \tilde{X}^{\text{Gauss}}U^\dagger$  with  $\tilde{X}$  being an  $N \times N$  Gaussian random matrix. Then by the definition of  $G$  in (2.28),

$$G(T, X^{\text{Gauss}}, z, w) \stackrel{d}{=} \begin{pmatrix} U & 0 \\ 0 & U \end{pmatrix} G(D, \tilde{X}^{\text{Gauss}}, z, w) \begin{pmatrix} U^\dagger & 0 \\ 0 & U^\dagger \end{pmatrix}. \quad (2.50)$$

Since the anisotropic local law holds for  $G(D, \tilde{X}^{\text{Gauss}}, z, w)$  by Theorem 2.18, we get immediately the anisotropic local law for  $G(T, X^{\text{Gauss}}, z, w)$ . The next theorem states that the anisotropic local law holds for general  $TX$  provided that the anisotropic local law holds for  $TX^{\text{Gauss}}$ . —

**Theorem 2.19** (Anisotropic local law when  $N \leq M$ ). Fix  $\tau \leq ||z|^2 - 1| \leq \tau^{-1}$ . Suppose Assumption 2.1 holds and  $N \leq M$ . Let  $T = U\bar{D}V$  be a singular decomposition of  $T$ , where  $\bar{D} = (D, 0)$  with  $D = \text{diag}(d_1, d_2, \dots, d_N)$ . Let  $\mathbf{S}$  be a regular domain. Then the anisotropic local law and averaged local law hold with parameters  $(T, X, z, \mathbf{S} \cap \hat{\mathbf{D}})$ . If in addition (2.23) holds, then the anisotropic local law and averaged local law hold with parameters  $(T, X, z, \mathbf{S})$ .

Finally we turn to the  $N > M$  case. Suppose  $T = U\bar{D}V$  is a singular decomposition of  $T$ , where  $U$  is an  $N \times N$  unitary matrix,  $V$  is an  $M \times M$  unitary matrix and  $\bar{D} = \begin{pmatrix} D \\ 0 \end{pmatrix}$  is an  $N \times M$

matrix such that  $D = \text{diag}(d_1, d_2, \dots, d_M)$ . Let  $U = (U_1, U_2)$ , where  $U_1$  has size  $N \times M$  and  $U_2$  has size  $N \times (N - M)$ . Following Girko's idea of Hermitization [20], to prove the local circular law in Theorem 2.6 when  $N > M$ , it suffices to study  $\det(TX - z)$  (see (2.52) below), for which we have

$$\det(TX - z) = \det \begin{pmatrix} DVXU_1 - z & DVXU_2 \\ 0 & -z \end{pmatrix} = \det(V^T D^T U_1^T X^T - z)(-z)^{N-M}. \quad (2.51)$$

Comparing with (2.49), we see that this case is reduced to the  $N \leq M$  case, with the only difference being that the extra  $(-z)^{N-M}$  term corresponds to the  $N - M$  zero eigenvalues of  $TX$ . Thus we make the following claim.

**Claim 2.20.** *The  $N < M$  case of Theorem 2.6 implies the  $N > M$  case of Theorem 2.6.*

## 2.4 Proof of Theorem 2.6

By Claim 2.20, it suffices to assume  $N \leq M$ . Our main tool will be Theorem 2.19. A major part of the proof follows from [7, Section 5]. The following lemma collects basic properties of stochastic domination  $<$ , which will be used tacitly during the proof and throughout this paper.

**Lemma 2.21** (Lemma 3.2 in [4]). *(i) Suppose that  $\xi(u, v) < \zeta(u, v)$  uniformly in  $u \in U$  and  $v \in V$ . If  $|V| \leq N^C$  for some constant  $C$ , then*

$$\sum_{v \in V} \xi(u, v) < \sum_{v \in V} \zeta(u, v)$$

*uniformly in  $u$ .*

*(ii) If  $\xi_1(u) < \zeta_1(u)$  uniformly in  $u \in U$  and  $\xi_2(u) < \zeta_2(u)$  uniformly in  $u \in U$ , then*

$$\xi_1(u)\xi_2(u) < \zeta_1(u)\zeta_2(u)$$

*uniformly in  $u \in U$ .*

*(iii) Suppose that  $\Psi(u) \geq N^{-C}$  is deterministic and  $\xi(u)$  is a nonnegative random variable such that  $E\xi(u)^2 \leq N^C$  for all  $u$ . Then if  $\xi(u) < \Psi(u)$  uniformly in  $u$ , we have*

$$\mathbb{E}\xi(u) < \Psi(u)$$

*uniformly in  $u$ .*

The Girko's Hermitization technique [20] can be reformulated as the following (see e.g. [22]): for any smooth function  $g$ ,

$$\begin{aligned} \frac{1}{N} \sum_{j=1}^N g(\mu_j) &= \frac{1}{4\pi N} \int \Delta g(z) \sum_{j=1}^N \log(\mu_j - z)(\bar{\mu}_j - \bar{z}) dA(z) \\ &= \frac{1}{4\pi N} \int \Delta g(z) \log |\det(Y(z)Y^\dagger(z))| dA(z) = \frac{1}{4\pi N} \int \Delta g(z) \sum_{j=1}^N \log \lambda_j(z) dA(z), \end{aligned} \quad (2.52)$$

where  $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N$  are the ordered eigenvalues of  $Y(z)Y^\dagger(z)$ . For  $g = F_{z_0, a}$ , we use the new variable  $\xi = N^a(z - z_0)$  to write the above equation as

$$\frac{1}{N} \sum_{j=1}^N F_{z_0, a}(\mu_j) = \frac{N^{-1+2a}}{4\pi} \int (\Delta F)(\xi) \sum_{j=1}^N \log \lambda_j(z) dA(\xi). \quad (2.53)$$

Define the classical location  $\gamma_j(z)$  of the  $j$ -th eigenvalue of  $Y(z)Y^\dagger(z)$  by

$$\int_0^{\gamma_j(z)} \rho_{2c}(x) dx = \frac{j}{N}, \quad 1 \leq j \leq N. \quad (2.54)$$

By Proposition 2.14, we have that for any  $\delta > 0$

$$\left| \sum_{j=1}^N \log \gamma_j(z) - N \int_0^\infty (\log x) \rho_{2c}(x, z) dx \right| \leq \sum_{j=1}^N N \int_{\gamma_{j-1}(z)}^{\gamma_j(z)} |\log \gamma_j(z) - \log x| \rho_{2c}(x, z) dx \leq N^\delta \quad (2.55)$$

for large enough  $N$ . Suppose we have the bound

$$\left| \sum_j \log \lambda_j - \sum_j \log \gamma_j \right| < N^b. \quad (2.56)$$

Plugging (2.55) and (2.56) into (2.53), we get

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N F_{z_0}(\mu_j) &= \frac{N^{2a}}{4\pi} \int (\Delta F)(\xi) \int_0^\infty (\log x) \rho_{2c}(x, z) dx dA(\xi) + O_{<}(N^{-1+b+2a} \|\Delta F\|_{L_1}) \\ &= \frac{1}{4\pi} \int F(\xi) \int_0^\infty (\log x) \Delta_z \rho_{2c}(x, z) dx dA(\xi) + O_{<}(N^{-1+b+2a} \|\Delta F\|_{L_1}). \end{aligned}$$

Thus we obtain (2.21) if we can prove (2.56) for  $b = 1/2$ , and we obtain (2.24) if we can prove (2.56) for  $b = 0$  when  $1 + \tau \leq |z_0|^2 \leq 1 + \tau^{-1}$  or the assumption (2.23) holds.

We need the following lemma which is a consequence of Theorem 2.19. Recall (2.16) and (2.20), the number of components  $L$  has order 1 and each component  $[e_{2k}, e_{2k-1}]$  contains order  $N$  of  $\gamma_j$ 's. We define the classical number of eigenvalues to the left of the edge  $e_k$ ,  $1 \leq k \leq 2L$ , as

$$N_k := \left\lfloor N \int_0^{e_k} \rho_{2c}(x) dx \right\rfloor. \quad (2.57)$$

Note that  $N_{2L} = 0$ ,  $N_1 = N$  and  $N_{2k+1} = N_{2k}$ ,  $1 \leq k \leq L-1$ .

**Lemma 2.22** (Singular value rigidity). *Fix a small  $\epsilon > 0$ .*

(i) *If the averaged local law holds with parameters  $(T, X, z, \mathbf{D}(\zeta, N) \cap \widehat{\mathbf{D}}(\zeta, N))$  for arbitrarily small  $\zeta$ , then the following estimates hold. For any  $e_{2k} > 0$  and  $N_{2k} + N^{1/2+\epsilon} \leq j \leq N_{2k-1} - N^{1/2+\epsilon}$ ,*

$$\frac{|\lambda_j - \gamma_j|}{\gamma_j} < \left( \min \left\{ \frac{j - N_{2k}}{N}, \frac{N_{2k-1} - j}{N} \right\} \right)^{-1/3} N^{-1/2}. \quad (2.58)$$

*In the case  $|z|^2 \leq 1 - \tau$  with  $e_{2L} = 0$ , we have for any  $N_{2L} + N^{1/2+\epsilon} \leq j \leq N_{2L-1} - N^{1/2+\epsilon}$ ,*

$$\frac{|\lambda_j - \gamma_j|}{\gamma_j} < j^{-1} \left( \frac{N_{2L-1} - j}{N} \right)^{-1/3} N^{1/2}. \quad (2.59)$$

*Moreover, if  $1 + \tau \leq |z|^2 \leq 1 + \tau^{-1}$ , then for any fixed  $0 < c < e_{2L}$ ,*

$$\#\{j : 0 < \lambda_j < c\} < 1. \quad (2.60)$$

(ii) If the averaged local law holds with parameters  $(T, X, z, \mathbf{D}(\zeta, N))$  for arbitrarily small  $\zeta$ , then the following estimates hold. For any  $e_{2k} > 0$  and  $N_{2k} + N^\epsilon \leq j \leq N_{2k-1} - N^\epsilon$ ,

$$\frac{|\lambda_j - \gamma_j|}{\gamma_j} < \left( \min \left\{ \frac{j - N_{2k}}{N}, \frac{N_{2k-1} - j}{N} \right\} \right)^{-1/3} N^{-1}. \quad (2.61)$$

In the case  $|z|^2 \leq 1 - \tau$  with  $e_{2L} = 0$ , we have for any  $N_{2L} + N^\epsilon \leq j \leq N_{2L-1} - N^\epsilon$ ,

$$\frac{|\lambda_j - \gamma_j|}{\gamma_j} < j^{-1} \left( \frac{N_{2L-1} - j}{N} \right)^{-1/3}. \quad (2.62)$$

*Proof.* The proof is similar to the proof of [7, Lemma 5.1]. See also [4, Theorem 2.10] or [14, Theorem 7.6]  $\square$

Using (2.58) and (2.59), we get that

$$\sum_{N_{2k} + N^{1/2+\epsilon} \leq j \leq N_{2k-1} - N^{1/2+\epsilon}} |\log \lambda_j - \log \gamma_j| < \sum_{N_{2k} + N^{1/2+\epsilon} \leq j \leq N_{2k-1} - N^{1/2+\epsilon}} \frac{|\lambda_j - \gamma_j|}{\gamma_j} < N^{1/2}. \quad (2.63)$$

Through a standard large deviation estimate, we have the following bound (see e.g. [9, 27, 33]),

$$\mathbb{P}(\|X\| > t) \leq e^{-c_0 t^2 N} \quad \text{for } t > C_0, \quad (2.64)$$

where  $c_0, C_0 > 0$  are constants. Thus we have

$$\lambda_j \leq \|Y\|^2 \leq (\|T\| \|X\| + |z|)^2 < 1, \quad 1 \leq j \leq N. \quad (2.65)$$

Together with Lemma 2.23 concerning the smallest singular value of  $TX - z$ , we get

$$\sum_{k=1}^{2L} \sum_{|j - e_k| < N^{1/2+\epsilon}} |\log \lambda_j| < N^{1/2+\epsilon}. \quad (2.66)$$

Since  $|\log \gamma_j| < 1$  by Proposition 2.14, we conclude

$$\sum_{k=1}^{2L} \sum_{|j - e_k| < N^{1/2+\epsilon}} |\log \lambda_j - \log \gamma_j| < N^{1/2+\epsilon}. \quad (2.67)$$

Combining (2.63)-(2.67), we get for any  $\epsilon > 0$ ,

$$\sum_{1 \leq j \leq N} |\log \lambda_j - \log \gamma_j| < N^{1/2+\epsilon} \quad (2.68)$$

for large enough  $N$ . This implies (2.56) for  $b = 1/2$ . If in addition the assumption (2.23) holds, the averaged local law holds with parameters  $(T, X, z, \mathbf{D}(\zeta, N))$  for arbitrarily small  $\zeta$  by Theorem 2.19. Then we can prove (2.56) for  $b = 0$  using the better bounds (2.61) and (2.62).

Finally we prove that when  $|z_0|^2 \geq 1 + \tau$ , with the bounds (2.58) we can still prove the estimate (2.56) for  $b = 0$ . By the averaged local law and the definition of  $\gamma_j$  in (2.54), we have

$$\left| \sum_{j=1}^N \frac{1}{\lambda_j - i\eta} - \sum_{j=1}^N \frac{1}{\gamma_j - i\eta} \right| < \frac{1}{\eta}, \quad (2.69)$$



uniformly in  $N^{-1/2+\epsilon} \leq \eta \leq N^{1/2}$ . Taking integral of (2.69) over  $\eta$  from  $N^{-1/2+\epsilon}$  to  $N^{1/2}$ , we get

$$\left| \sum_{j=1}^N \log \left( \frac{\lambda_j - iN^{-1/2+\epsilon}}{\gamma_j - iN^{-1/2+\epsilon}} \right) - \sum_{j=1}^N \log \left( \frac{\lambda_j - iN^{1/2}}{\gamma_j - iN^{1/2}} \right) \right| < 1. \quad (2.70)$$

Then we use (2.58) and the bound (2.65) to estimate that

$$\left| \sum_{j=1}^N \log \left( \frac{\lambda_j - iN^{1/2}}{\gamma_j - iN^{1/2}} \right) \right| < \sum_{j=1}^N \left| (\lambda_j - \gamma_j) N^{-1/2} \right| < N^\epsilon.$$

Thus we conclude

$$\left| \sum_{j=1}^N \log \left( \frac{\lambda_j - iN^{-1/2+\epsilon}}{\gamma_j - iN^{-1/2+\epsilon}} \right) \right| < N^\epsilon. \quad (2.71)$$

Using  $\gamma_j \sim 1$ , (2.60) and (2.73), we get

$$\begin{aligned} \left| \sum_{j=1}^N \log \left( \frac{\lambda_j - iN^{-1/2+\epsilon}}{\gamma_j - iN^{-1/2+\epsilon}} \right) - \sum_{j=1}^N \log \frac{\lambda_j}{\gamma_j} \right| &< 1 + \left| \sum_{\lambda_j \geq c} \log \left( \frac{\lambda_j - iN^{-1/2+\epsilon}}{\gamma_j - iN^{-1/2+\epsilon}} \right) - \sum_{\lambda_j \geq c} \log \frac{\lambda_j}{\gamma_j} \right| \\ &< 1 + \sum_{\lambda_j \geq c} \left| (\lambda_j - \gamma_j) N^{-1/2+\epsilon} \right| < N^{2\epsilon}. \end{aligned} \quad (2.72)$$

Combing (2.71) and (2.72), we conclude (2.56) for  $b = 0$ .

**Lemma 2.23** (Lower bound on the smallest singular value). *If  $N < M$  and the entries of  $X$  have a density bounded by  $N^{C_3}$  for some  $C_3 > 0$ , then*

$$|\log \lambda_1(z)| < 1 \quad (2.73)$$

*holds uniformly for  $z$  in any fixed compact set. If  $N = M$ , the bounded density condition is not necessary.*

*Proof.* To prove (2.73), we need to prove that

$$\mathbb{P} \left( \lambda_1(z) \leq e^{-N^\epsilon} \right) \leq N^{-C} \quad (2.74)$$

for any  $\epsilon, C > 0$ . In the case  $N = M$  without the bounded density assumption, we have  $\lambda_1(z) \geq \tau \lambda'_1(z)$ , where  $\lambda'_1(z)$  is the smallest singular values of  $X - T^{-1}z$ . Following [32] or [36, Theorem 2.1], we have  $|\log \lambda'_1(z)| < 1$ , which further proves (2.73).

Now we turn to the case  $N < M$  with the bounded density assumption. By (2.49) we have that

$$TX - z = UD(V_1X - D^{-1}U^{-1}z) =: UD\tilde{Y}(z).$$

Hence it suffices to control the smallest singular value of  $\tilde{Y}(z)$ , call it  $\tilde{\lambda}_1(z)$ . Notice the columns  $\tilde{Y}_1, \dots, \tilde{Y}_N$  of  $\tilde{Y}(z)$  are independent vectors. From the variational characterization

$$\tilde{\lambda}_1(z) = \min_{|u|=1} \|\tilde{Y}(z)u\|^2,$$

we can easily get

$$\tilde{\lambda}_1(z)^{1/2} \geq N^{-1/2} \min_{1 \leq k \leq N} \text{dist} \left( \tilde{Y}_k, \text{span}\{\tilde{Y}_l, l \neq k\} \right) = N^{-1/2} \min_{1 \leq k \leq N} |\langle \tilde{Y}_k, u_k \rangle|, \quad (2.75)$$

where  $u_k$  is the unit normal vector of  $\text{span}\{\tilde{Y}_l, l \neq k\}$  and hence is independent of  $\tilde{Y}_k$ . By conditioning on  $u_k$ , we get immediately

$$\mathbb{P}(\tilde{\lambda}_1(z) \leq N^{-C_0}) \leq CN^{-C_0/2+C_3+3/2}, \quad (2.76)$$

which is a much stronger result than (2.74). Here we have used Theorem 1.2 of [34] to conclude that  $\langle \tilde{Y}_k, u_k \rangle$  for fixed  $u_k$  has density bounded by  $CN^{C_3}$ .  $\square$

## 2.5 Outline of the paper

The rest of this paper is devoted to the proof of Theorems 2.18 and 2.19. In Section 3, we collect the basics tools that we shall use throughout the proof. In Section 4, we perform step (A) of the proof by proving the entrywise local law and averaged local law in Theorem 2.18 under the assumption that  $T$  is diagonal. We first prove a weak version of the entrywise local law in Sections 4.1-4.3, and then improve the weak law to the strong entrywise local law and averaged local law in Sections 4.4-4.5. In Section 5, we perform step (B) of the proof by proving the anisotropic local law in Theorem 2.18 using the entrywise local law proved in Section 4. Finally in Section 6 we finish the step (C) of the proof, where using Theorem 2.18, we prove Theorem 2.19 with a self-consistent comparison method.

The first part of Appendix A establishes the basic properties of  $\rho_{1,2c}$  stated in Lemma 2.3 and Proposition 2.14. In Sections A.2 and A.3, we establish some key estimates on  $m_{1,2c}$  and the stability of the self-consistent equation (2.11) on regular domains.

## 3 Basic tools

In this preliminary section, we collect various identities and estimates that we shall use throughout the following.

**Definition 3.1** (Minors). For  $J \subset \mathcal{I}$ , we define the minor  $H^{(J)} := \{H_{st} : s, t \in \mathcal{I} \setminus J\}$ , and correspondingly  $G^{(J)} := (H^{(J)})^{-1}$ . Let  $[J] := \{s \in \mathcal{I} : s \in J \text{ or } \bar{s} \in J\}$ . We also denote  $H^{[J]} := \{H_{st} : s, t \in \mathcal{I} \setminus [J]\}$  and  $G^{[J]} := (H^{[J]})^{-1}$ . We abbreviate  $(\{s\}) \equiv (s)$ ,  $(\{s, t\}) \equiv (st)$ ,  $[\{s\}] \equiv [s]$  and  $[\{s, t\}] \equiv [st]$ .

Notice that by the definition, we have  $H_{st}^{(J)} = 0$  and  $G_{st}^{(J)} = 0$  if  $s \in J$  or  $t \in J$ .

**Lemma 3.2.** (Resolvent identities).

(i) For  $i \in \mathcal{I}_1$  and  $\mu \in \mathcal{I}_2$ , we have

$$\frac{1}{G_{ii}} = -w - w \left( Y G^{(i)} Y^\dagger \right)_{ii}, \quad \frac{1}{G_{\mu\mu}} = -w - w \left( Y^\dagger G^{(\mu)} Y \right)_{\mu\mu}. \quad (3.1)$$

For  $i \neq j \in \mathcal{I}_1$  and  $\mu \neq \nu \in \mathcal{I}_2$ , we have

$$G_{ij} = w G_{ii} G_{jj}^{(i)} \left( Y G^{(ij)} Y^\dagger \right)_{ij}, \quad G_{\mu\nu} = w G_{\mu\mu} G_{\nu\nu}^{(\mu)} \left( Y^\dagger G^{(\mu\nu)} Y \right)_{\mu\nu}. \quad (3.2)$$

(ii) For  $i \in \mathcal{I}_1$  and  $\mu \in \mathcal{I}_2$ , we have

$$G_{i\mu} = G_{ii} G_{\mu\mu}^{(i)} \left( -w^{1/2} Y_{i\mu} + w \left( Y G^{(i\mu)} Y \right)_{i\mu} \right), \quad (3.3)$$

$$G_{\mu i} = G_{\mu\mu} G_{ii}^{(\mu)} \left( -w^{1/2} Y_{\mu i}^\dagger + w \left( Y^\dagger G^{(\mu i)} Y^\dagger \right)_{\mu i} \right). \quad (3.4)$$

(iii) For  $r \in \mathcal{I}$  and  $s, t \in \mathcal{I} \setminus \{r\}$ ,

$$G_{st}^{(r)} = G_{st} - \frac{G_{sr}G_{rt}}{G_{rr}}, \quad \frac{1}{G_{ss}} = \frac{1}{G_{ss}^{(r)}} - \frac{G_{sr}G_{rs}}{G_{ss}G_{ss}^{(r)}G_{rr}}. \quad (3.5)$$

(iv) All of the above identities hold for  $G^{(J)}$  instead of  $G$  for  $J \subset \mathcal{I}$ .

*Proof.* All these identities can be proved using Schur's complement formula. They have been previously derived and summarized e.g. in [14, 15, 17].  $\square$

**Lemma 3.3.** (Resolvent identities for  $G_{[ij]}$  groups).

(i) For  $i \in \mathcal{I}_1$ , we have

$$G_{[ii]}^{-1} = H_{[ii]} - \sum_{k, l \neq i} H_{[ik]} G_{[kl]}^{[i]} H_{[li]}. \quad (3.6)$$

For  $i \neq j \in \mathcal{I}_1$ , we have

$$G_{[ij]} = -G_{[ii]} \sum_{k \neq i} H_{[ik]} G_{[kj]}^{[i]} = - \sum_{k \neq j} G_{[ik]}^{[j]} H_{[kj]} G_{[jj]} \quad (3.7)$$

$$= -G_{[ii]} H_{[ij]} G_{[jj]}^{[i]} + G_{[ii]} \sum_{k, l \notin \{i, j\}} H_{[ik]} G_{[kl]}^{[ij]} H_{[lj]} G_{[jj]}^{[i]}. \quad (3.8)$$

(ii) For  $k \in \mathcal{I}_1$  and  $i, j \in \mathcal{I}_1 \setminus \{k\}$ ,

$$G_{[ij]}^{[k]} = G_{[ij]} - G_{[ik]} G_{[kk]}^{-1} G_{[kj]}, \quad (3.9)$$

and

$$G_{[ii]}^{-1} = \left( G_{[ii]}^{[k]} \right)^{-1} - G_{[ii]}^{-1} G_{[ik]} G_{[kk]}^{-1} G_{[ki]} \left( G_{[ii]}^{[k]} \right)^{-1}. \quad (3.10)$$

(iii) All of the above identities hold for  $G^{[J]}$  instead of  $G$  for  $J \subset \mathcal{I}$ .

*Proof.* These identities can be proved using Schur's complement formula. The details are left to the reader.  $\square$

Next we introduce the spectral decomposition of  $G$ . Let

$$Y = \sum_{k=1}^N \sqrt{\lambda_k} \xi_k \zeta_k^\dagger$$

be the singular decomposition of  $Y$ , where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_N \geq 0$  and  $\{\xi_k\}_{k=1}^N$  and  $\{\zeta_k\}_{k=1}^N$  are orthonormal bases of  $\mathbb{C}^{\mathcal{I}_1}$  and  $\mathbb{C}^{\mathcal{I}_2}$  respectively. Then by (2.30), we have

$$G(w) = \sum_{k=1}^N \frac{1}{\lambda_k - w} \begin{pmatrix} \xi_k \xi_k^\dagger & w^{-1/2} \sqrt{\lambda_k} \xi_k \zeta_k^\dagger \\ w^{-1/2} \sqrt{\lambda_k} \zeta_k \xi_k^\dagger & \zeta_k \zeta_k^\dagger \end{pmatrix}. \quad (3.11)$$

**Definition 3.4** (Generalized entries). For  $\mathbf{v}, \mathbf{w} \in \mathbb{C}^{\mathcal{I}}$ ,  $s \in \mathcal{I}$  and an  $\mathcal{I} \times \mathcal{I}$  matrix  $A$ , we shall denote

$$A_{\mathbf{v}\mathbf{w}} := \langle \mathbf{v}, A\mathbf{w} \rangle, \quad A_{\mathbf{v}s} := \langle \mathbf{v}, A\mathbf{e}_s \rangle, \quad A_{s\mathbf{w}} := \langle \mathbf{e}_s, A\mathbf{w} \rangle, \quad (3.12)$$

where  $\mathbf{e}_s$  is the standard unit vector.

Given vectors  $\mathbf{v} \in \mathbb{C}^{\mathcal{I}_1}$  and  $\mathbf{w} \in \mathbb{C}^{\mathcal{I}_2}$ , we always identify them with their natural embeddings  $\begin{pmatrix} \mathbf{v} \\ 0 \end{pmatrix}$  and  $\begin{pmatrix} 0 \\ \mathbf{w} \end{pmatrix}$  in  $\mathbb{C}^{\mathcal{I}}$ . The exact meanings will be clear from the context.

**Lemma 3.5.** Fix  $\tau > 0$ . The following estimates hold uniformly for any  $w \in \mathbf{D}(\tau, N)$ . We have

$$\|G\| \leq C\eta^{-1}, \quad \|\partial_w G\| \leq C\eta^{-2}. \quad (3.13)$$

Let  $\mathbf{v} \in \mathbb{C}^{\mathcal{I}_1}$  and  $\mathbf{w} \in \mathbb{C}^{\mathcal{I}_2}$ , we have the bounds

$$\sum_{\mu \in \mathcal{I}_2} |G_{\mathbf{w}\mu}|^2 = \sum_{\mu \in \mathcal{I}_2} |G_{\mu\mathbf{w}}|^2 = \frac{\text{Im } G_{\mathbf{w}\mathbf{w}}}{\eta}, \quad (3.14)$$

$$\sum_{i \in \mathcal{I}_1} |G_{\mathbf{v}i}|^2 = \sum_{i \in \mathcal{I}_1} |G_{i\mathbf{v}}|^2 = \frac{\text{Im } G_{\mathbf{v}\mathbf{v}}}{\eta}, \quad (3.15)$$

$$\sum_{i \in \mathcal{I}_1} |G_{\mathbf{w}i}|^2 = \sum_{i \in \mathcal{I}_1} |G_{i\mathbf{w}}|^2 = |w|^{-1} G_{\mathbf{w}\mathbf{w}} + \bar{w} |w|^{-1} \frac{\text{Im } G_{\mathbf{w}\mathbf{w}}}{\eta}, \quad (3.16)$$

$$\sum_{\mu \in \mathcal{I}_2} |G_{\mathbf{v}\mu}|^2 = \sum_{\mu \in \mathcal{I}_2} |G_{\mu\mathbf{v}}|^2 = |w|^{-1} G_{\mathbf{v}\mathbf{v}} + \bar{w} |w|^{-1} \frac{\text{Im } G_{\mathbf{v}\mathbf{v}}}{\eta}. \quad (3.17)$$

All of the above estimates remain true for  $G^{(J)}$  instead of  $G$  for  $J \subset \mathcal{I}$ .

*Proof.* The estimates in (3.13) follow from (3.11). For any unit vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{C}^{\mathcal{I}_1}$ , we have

$$|\langle \mathbf{x}, G\mathbf{y} \rangle| \leq \sum_{k=1}^N \frac{|\langle \mathbf{x}, \xi_k \rangle| |\langle \xi_k^\dagger, \mathbf{y} \rangle|}{|\lambda_k - w|} \leq \frac{1}{\eta} \left[ \sum_{k=1}^N |\langle \mathbf{x}, \xi_k \rangle|^2 \right]^{1/2} \left[ \sum_{k=1}^N |\langle \xi_k^\dagger, \mathbf{y} \rangle|^2 \right]^{1/2} = \frac{1}{\eta}.$$

For any unit vectors  $\mathbf{x} \in \mathbb{C}^{\mathcal{I}_1}$  and  $\mathbf{y} \in \mathbb{C}^{\mathcal{I}_2}$ , we have

$$|\langle \mathbf{x}, G\mathbf{y} \rangle| \leq |w|^{-1/2} \sum_{k=1}^N \frac{\sqrt{\lambda_k} |\langle \mathbf{x}, \xi_k \rangle| |\langle \zeta_k^\dagger, \mathbf{y} \rangle|}{|\lambda_k - w|} \leq \sum_{k=1}^N \frac{1}{2\eta} \left( |\langle \mathbf{x}, \xi_k \rangle|^2 + |\langle \zeta_k^\dagger, \mathbf{y} \rangle|^2 \right) = \frac{1}{\eta},$$

where we have used that for  $w = E + i\eta$ ,  $|w|^{-1/2} \sqrt{\lambda_k}/|\lambda_k - w| \leq \eta^{-1}$ . For the other two blocks of  $G$ , we can prove similar estimates. This implies (3.13). It is trivial to generalize the proof to  $\partial_w G$ , where  $\eta^{-2}$  comes from the  $(\lambda_k - w)^{-2}$  factor of  $\partial_w G$ . For (3.14), we observe that

$$\frac{\text{Im } G_{\mathbf{w}\mathbf{w}}}{\eta} = \frac{1}{\eta} \text{Im} \sum_{k=1}^N \frac{\langle \mathbf{w}, \zeta_k \rangle \langle \zeta_k^\dagger, \mathbf{w} \rangle}{\lambda_k - w} = \sum_{k=1}^N \frac{|\langle \mathbf{w}, \zeta_k \rangle|^2}{(\lambda_k - E)^2 + \eta^2}$$

and by (2.30)

$$\sum_{\mu \in \mathcal{I}_2} |G_{\mathbf{w}\mu}|^2 = \sum_{\mu \in \mathcal{I}_2} \langle \mathbf{w}, G_R e_\mu \rangle \langle e_\mu, G_R^\dagger \mathbf{w} \rangle = \left\langle \mathbf{w}, G_R G_R^\dagger \mathbf{w} \right\rangle = \sum_{k=1}^N \frac{|\langle \mathbf{w}, \zeta_k \rangle|^2}{(\lambda_k - E)^2 + \eta^2}. \quad (3.18)$$

Similarly, we can prove the identity for  $\sum_{\mu \in \mathcal{I}_2} |G_{\mu\mathbf{w}}|^2$  and (3.15). For identity (3.16), first we can prove

$\sum_{i \in \mathcal{I}_1} |G_{\mathbf{w}i}|^2 = \sum_{i \in \mathcal{I}_1} |G_{i\mathbf{w}}|^2$  using (3.11). Then we use (2.30) and (3.18) to get

$$\sum_{i \in \mathcal{I}_1} |G_{\mathbf{w}i}|^2 = |w|^{-1} \left( G_R Y^\dagger Y G_R^\dagger \right)_{\mathbf{w}\mathbf{w}} = |w|^{-1} \left[ G_R (Y^\dagger Y - \bar{w}) G_R^\dagger \right]_{\mathbf{w}\mathbf{w}} + \bar{w} |w|^{-1} \left( G_R G_R^\dagger \right)_{\mathbf{w}\mathbf{w}}$$

$$= |w|^{-1} G_{\mathbf{w}\mathbf{w}} + \bar{w} |w|^{-1} \left( G_R G_R^\dagger \right)_{\mathbf{w}\mathbf{w}} = |w|^{-1} G_{\mathbf{w}\mathbf{w}} + \bar{w} |w|^{-1} \frac{\text{Im } G_{\mathbf{w}\mathbf{w}}}{\eta}. \quad (3.19)$$

Identity (3.17) can be proved in a similar way.  $\square$

The following Lemma give useful large deviation bounds. See Theorem B.1 and Lemmas B.2-B.4 in [13] for the proof. See also Theorem C.1 of [14].

**Lemma 3.6.** *(Large deviation bounds) Let  $(X_i^{(N)})$ ,  $(Y_i^{(N)})$  be independent families of random variables and  $(a_{ij}^{(N)})$ ,  $(b_i^{(N)})$  be deterministic. Suppose all entries  $X_i^{(N)}$  and  $Y_i^{(N)}$  are independent and satisfies (2.2) and (2.3). Then we have the following bounds:*

$$\sum_i b_i X_i < \frac{\left( \sum_i |b_i|^2 \right)^{1/2}}{\sqrt{N}}, \quad \sum_{i,j} a_{ij} X_i Y_j < \frac{\left( \sum_{i,j} |a_{ij}|^2 \right)^{1/2}}{N}, \quad \sum_{i \neq j} a_{ij} X_i Y_j < \frac{\left( \sum_{i \neq j} |a_{ij}|^2 \right)^{1/2}}{N}. \quad (3.20)$$

If the coefficients  $(a_{ij}^{(N)})$  and  $(b_i^{(N)})$  depend on some parameter  $u$ , then all of the above estimates are uniform in  $u$ .

We have stated some basic properties of  $\rho_{1,2c}$  and  $m_{1,2c}$  in Lemma 2.3 and Proposition 2.14. Now we collect more estimates for  $m_{1,2c}$  that will be used in the proof. The next lemma is proved in Appendix A.2. For  $w = E + i\eta \in \mathbf{D}$ , we define the distance to the spectral edge through

$$\kappa \equiv \kappa(E) := \min_{1 \leq k \leq 2L, e_k > 0} |E - e_k|. \quad (3.21)$$

Notice in the  $|z| < 1$  case, we do not take into consideration the edge at  $e_{2L} = 0$ .

**Lemma 3.7.** *Fix  $\tau > 0$  and suppose  $\tau \leq ||z|^2 - 1| \leq \tau^{-1}$ . We denote  $w = E + i\eta$ .*

*Case 1 Fix  $\tau' > 0$ . Suppose the bulk component  $[e_{2k}, e_{2k-1}]$  is regular in the sense of Definition 2.4. Then for  $w \in \mathbf{D}_k^b(\zeta, \tau', N)$ , we have*

$$|1 + m_{1c}| \sim \text{Im } m_{1c} \sim 1, \quad |m_{2c}| \sim \text{Im } m_{2c} \sim 1. \quad (3.22)$$

*Case 2 Fix  $\tau' > 0$ . Then for  $w \in \mathbf{D}^o(\zeta, \tau', N)$ , we have*

$$\text{Im } m_{1,2c} \sim \eta, \quad |1 + m_{1c}| \sim 1, \quad |m_{2c}| \sim 1. \quad (3.23)$$

*Case 3 Suppose  $e_k \neq 0$  is a regular edge. Then for  $w \in \mathbf{D}_k^e(\zeta, \tau', N)$ , if  $\tau' > 0$  is small enough,*

$$\text{Im } m_{1,2c} \sim \begin{cases} \sqrt{\kappa + \eta} & \text{if } E \in \text{supp } \rho_{1,2c}, \\ \eta/\sqrt{\kappa + \eta} & \text{if } E \notin \text{supp } \rho_{1,2c} \end{cases}, \quad |1 + m_{1c}| \sim 1, \quad |m_{2c}| \sim 1. \quad (3.24)$$

*Case 4 Suppose  $|z|^2 \leq 1 - \tau$ . We take  $e_{2L} = 0$  and  $\tau' > 0$  to be small enough. Then for  $w \in \mathbf{D}_{2L}^e(\zeta, \tau', N)$ , if  $\text{Im } w \geq \tau'$ , we have*

$$|1 + m_{1c}| \sim \text{Im } m_{1c} \sim 1, \quad |m_{2c}| \sim \text{Im } m_{2c} \sim 1; \quad (3.25)$$

*if  $|w| \leq 2\tau'$ , we have*

$$m_{1c} = i \frac{\sqrt{t}}{\sqrt{w}} + O(1), \quad m_{2c} = \frac{i\sqrt{t}}{\sqrt{w}(t + |z|^2)} + O(1), \quad (3.26)$$

*for some constant  $t > 0$ , and*

$$\text{Im } m_{1,2c} \sim |w|^{-1/2}. \quad (3.27)$$

Case 5 For  $w \in \mathbf{D}_L(\zeta)$ , we have

$$|m_{1c}| \sim \operatorname{Im} m_{1c} \sim \frac{1}{\eta}, \quad |m_{2c}| \sim \operatorname{Im} m_{2c} \sim \frac{1}{\eta}. \quad (3.28)$$

In Cases 1-4, we have

$$\left| w(1 + s_i m_{2c})(1 + m_{1c}) - |z|^2 \right| \geq c, \quad (3.29)$$

where  $c > 0$  is some constant that may depend on  $\tau$  and  $\tau'$ . In Case 5, we have

$$\left| w(1 + s_i m_{2c})(1 + m_{1c}) - |z|^2 \right| \geq \eta, \quad (3.30)$$

Note that the uniform bounds (3.29) and (3.30) guarantee that the matrix entries of  $\Pi(w)$  remain bounded. We have the following Lemma, which is prove in Appendix A.2.

**Lemma 3.8.** *In Cases 1-4 of Lemma 3.7, we have*

$$\|\pi_{[i]c}\| \leq C|w|^{-1/2}, \quad \left\| (\pi_{[i]c})^{-1} \right\| \leq C|w|^{1/2}, \quad (3.31)$$

and in Case 5 of Lemma 3.7, we have

$$\|\pi_{[i]c}\| \leq C\eta^{-1}, \quad \left\| (\pi_{[i]c})^{-1} \right\| \leq C\eta. \quad (3.32)$$

For all the cases in Lemma 3.7,

$$\operatorname{Im} \Pi_{\mathbf{v}\mathbf{v}} \leq C\operatorname{Im}(m_{1c} + m_{2c}), \quad (3.33)$$

uniformly in  $w$  and any deterministic unit vector  $\mathbf{v} \in \mathbb{C}^{\mathcal{I}}$ .

The self-consistent equation (2.11) can be written as

$$\Upsilon(w, m_1) = 0, \quad (3.34)$$

where

$$\Upsilon(w, m_1) = m_1 + \frac{1}{N} \sum_{i=1}^n l_i s_i (1 + m_1) \left[ w \left( 1 + s_i \frac{1 + m_1}{-w(1 + m_1)^2 + |z|^2} \right) (1 + m_1) - |z|^2 \right]^{-1}. \quad (3.35)$$

The stability of (3.34) roughly says that if  $\Upsilon(w, m_1)$  is small and  $m_1(w') - m_{1c}(w')$  is small for  $w' := w + iN^{-10}$ , then  $m_1(w) - m_{1c}(w)$  is small. For an arbitrary  $w \in \mathbf{D}$ , we define the discrete set

$$L(w) := \{w\} \cup \{w' \in \mathbf{D} : \operatorname{Re} w' = \operatorname{Re} w, \operatorname{Im} w' \in [\operatorname{Im} w, 1] \cap (N^{-10}\mathbb{N})\}, \quad (3.36)$$

Thus, if  $\operatorname{Im} w \geq 1$  then  $L(w) = \{w\}$ , and if  $\operatorname{Im} w < 1$  then  $L(w)$  is a 1-dimensional lattice with spacing  $N^{-10}$  plus the point  $w$ . Obviously, we have  $|L(w)| \leq N^{10}$ .

**Definition 3.9** (Stability of (3.34)). *We say that (3.34) is stable on  $\mathbf{D}$  if the following holds. Suppose that  $N^{-2}|m_{1c}| \leq \delta(w) \leq (\log N)^{-1}|m_{1c}|$  for  $w \in \mathbf{D}$  and that  $\delta$  is Lipschitz continuous with Lipschitz constant  $\leq N^4$ . Suppose moreover that for each fixed  $E$ , the function  $\eta \mapsto \delta(E + i\eta)$  is*

non-increasing for  $\eta > 0$ . Suppose that  $u_1 : \mathbf{D} \rightarrow \mathbb{C}$  is the Stieltjes transform of a positive integrable function. Let  $w \in \mathbf{D}$  and suppose that for all  $w' \in L(w)$  we have

$$|\Upsilon(w, u_1)| \leq \delta(w). \quad (3.37)$$

Then

$$|u_1(w) - m_{1c}(w)| \leq \frac{C\delta}{\sqrt{\kappa + \eta + \delta}}, \quad (3.38)$$

for some constant  $C > 0$  independent of  $w$  and  $N$ .

We say that (3.34) is stable on  $\mathbf{D}_L$  if for  $0 \leq \delta(w) \leq (\log N)^{-1}|m_{1c}|$ , (3.37) implies

$$|u_1(w) - m_{1c}(w)| \leq C\delta, \quad (3.39)$$

for some constant  $C > 0$  independent of  $w$  and  $N$ .

This stability condition has previously appeared in [4, 7, 24]. In [24], for example, the stability condition was established under various regularity assumptions. In the following lemma, we establish the stability on each regular domain. The proof is presented in Appendix A.3. This lemma leaves the case  $|w|^{1/2} + |z|^2 = o(1)$  alone. We will handle this case in a different way in Section 4.5.

**Lemma 3.10.** Fix  $\tau > 0$  and let  $\tau' > 0$  be sufficiently small depending on  $\tau$ . Let  $\tau \leq ||z|^2 - 1| \leq \tau^{-1}$ .

*Case 1* Suppose the bulk component  $[e_{2k}, e_{2k-1}]$  is regular in the sense of Definition 2.4. Then (3.34) is stable on  $\mathbf{D}_k^b(\zeta, \tau', N)$  in the sense of Definition 3.9.

*Case 2* (3.34) is stable on  $\mathbf{D}^o(\zeta, \tau', N)$  in the sense of Definition 3.9.

*Case 3* Suppose  $e_k \neq 0$  is a regular edge in the sense of Definition 2.4. Then (3.34) is stable on  $\mathbf{D}_k^e(\zeta, \tau', N)$  in the sense of Definition 3.9.

*Case 4* Suppose  $|z|^2 \leq 1 - \tau$  and  $e_{2L} = 0$ . If  $|w|^{1/2} + |z|^2 \geq \epsilon$  for some constant  $\epsilon > 0$ , then (3.34) is stable on  $\mathbf{D}_{2L}^e(\zeta, \tau', N)$  in the sense of Definition 3.9.

*Case 5* (3.34) is stable on  $\mathbf{D}_L(\zeta)$  in the sense of Definition 3.9.

## 4 Entrywise local law when $T$ is diagonal

In this section we prove the entrywise local law and averaged local law in Theorem 2.18 when  $T$  is diagonal. The proof is similar to the previous proofs of entrywise local laws in e.g. [4, 5, 7, 24]. We basically follow the ideas in [7], and we will provide necessary details for the parts that are different from the previous proofs.

The main novel observation of this section is that the self-consistent equations (2.9) and (2.10) can be “derived” from the random matrix model by an application of Schur’s complement formula. It is helpful to give a heuristic argument here. We introduce the conditional expectation

$$\mathbb{E}_{[i]}[\cdot] := \mathbb{E}[\cdot \mid H^{[i]}],$$

i.e. the partial expectation in the randomness of the  $i$  and  $\bar{i}$ -th rows and columns of  $H$ . For the diagonal  $G_{[ii]}$  group, we ignore formally the random fluctuations in (3.6) to get that

$$G_{[ii]}^{-1} \approx \mathbb{E}_{[i]} H_{[ii]} - \sum_{k, l \neq i} \mathbb{E}_{[i]} \left( H_{[ik]} G_{[kl]}^{[i]} H_{[li]} \right) = \begin{pmatrix} -w & -w^{1/2}z \\ -w^{1/2}\bar{z} & -w \end{pmatrix} - \frac{w}{N} \sum_k \begin{pmatrix} |d_i|^2 G_{kk}^{[i]} & 0 \\ 0 & |d_k|^2 G_{kk}^{[i]} \end{pmatrix}$$

$$= \begin{pmatrix} -w & -w^{1/2}z \\ -w^{1/2}\bar{z} & -w \end{pmatrix} - w \begin{pmatrix} |d_i|^2 m_2 & 0 \\ 0 & m_1 \end{pmatrix}, \quad (4.1)$$

where we use the definition of  $m_1$  and  $m_2$  in (2.34). The 11 entry of (4.1) gives the equation

$$G_{ii} \approx \frac{-1 - m_1}{w(1 + |d_i|^2 m_2)(1 + m_1) - |z|^2}, \quad (4.2)$$

from which we get that

$$G_{ii} \left[ -w(1 + |d_i|^2 m_2) + \frac{|z|^2}{1 + m_1} \right] \approx 1.$$

Summing over  $i$  and using that  $N^{-1} \sum_i G_{ii} = N^{-1} \sum_\mu G_{\mu\mu} = m_2$ , the above equation becomes

$$-w(m_2 + m_1 m_2) + \frac{|z|^2 m_2}{1 + m_1} \approx 1,$$

which gives (2.9). Multiplying (4.2) with  $|d_i|^2$  and summing over  $i$ , we get the self-consistent equation (2.10). In this section we give a justification of these approximations.

Before we start the proof, we make the following remark. In this section we mainly focus on the domain  $\mathbf{D}$ . On the domain  $\mathbf{D}_L$ , the proofs are much simpler and we only describe them briefly. The parameter  $z$  can be either inside or outside of the unit circle. Recall Lemmas 3.7 and 3.10, the domain  $\mathbf{D}$  of  $w$  can be divided roughly into four cases:  $w$  near a *nonzero* regular edge,  $w \rightarrow 0$ ,  $w$  in the bulk, or  $w$  outside the spectrum. In this section we will only consider the case  $|z|^2 \leq 1 - \tau$  since it covers all four different behaviors. Notice in this case  $|m_{1,2c}(w)| \sim |w|^{-1/2}$  for  $w$  in any compact set of  $\mathbb{C}_+$  by Proposition 2.15. Also due to the remark above Lemma 3.10, in Sections 4.1-4.4, we assume  $|w|^{1/2} + |z|^2 \geq c$  for some  $c > 0$ . We will handle the  $|w|^{1/2} + |z|^2 = o(1)$  case in Section 4.5.

## 4.1 The self-consistent equations

To begin with, we prove the following weak version of the entrywise local law.

**Proposition 4.1** (Weak entrywise law). *Fix  $|z|^2 \leq 1 - \tau$  and a small constant  $c > 0$ . Suppose Assumption 2.1 holds,  $N = M$  and  $T \equiv D := \text{diag}(d_1, \dots, d_N)$ . Then for any regular domain  $\mathbf{S} \subset \mathbf{D}$ ,*

$$\max_{i,j \in \mathcal{I}_1} \left\| (G(w) - \Pi(w))_{[ij]} \right\| < \frac{1}{|w|^{1/2}} \left( \frac{|w|^{1/2}}{N\eta} \right)^{1/4} \quad (4.3)$$

for all  $w \in \mathbf{S}$  such that  $|w|^{1/2} + |z|^2 \geq c$ . For  $w \in \mathbf{D}_L$ , we have

$$\max_{i,j \in \mathcal{I}_1} \left\| (G(w) - \Pi(w))_{[ij]} \right\| < \frac{1}{\eta} \sqrt{\frac{1}{N}}. \quad (4.4)$$

For the purpose of proof, we define the following random control parameters.

**Definition 4.2** (Control parameters). *Suppose  $N = M$  and  $T \equiv D := \text{diag}(d_1, \dots, d_N)$ . We define*

$$\Lambda := \max_{i,j \in \mathcal{I}_1} \left\| (G - \Pi)_{[ij]} \right\|, \quad \Lambda_o := \max_{i \neq j \in \mathcal{I}_1} \left\| (G - \Pi)_{[ij]} \right\|. \quad (4.5)$$

For  $J \subseteq \mathcal{I}$ , define the averaged variables  $m_{1,2}^{(J)}$  ( $m_{1,2}^{[J]}$ ) by replacing  $G$  in (2.34) with  $G^{(J)}$  ( $G^{[J]}$ ), i.e.

$$m_1^{(J)} := \frac{1}{N} \sum_{i \notin J} |d_i|^2 G_{ii}^{(J)}, \quad m_2^{(J)} := \frac{1}{N} \sum_{\mu \notin J} G_{\mu\mu}^{(J)}. \quad (4.6)$$



The averaged error and the random control parameter are defined as

$$\theta := |m_1 - m_{1c}| + |m_2 - m_{2c}|, \quad \Psi_\theta := \sqrt{\frac{\operatorname{Im}(m_{1c} + m_{2c}) + \theta}{N\eta}} + \frac{1}{N\eta}. \quad (4.7)$$

*Remark:* By (2.4), we immediately get that

$$\tau \operatorname{Im} m_1^{(J)} \leq \operatorname{Im} m_2^{(J)} \leq \tau^{-1} \operatorname{Im} m_2^{(J)}, \quad (4.8)$$

and  $\theta = O(\Lambda)$ , since  $|m_1 - m_{1c}| \leq \tau^{-1}\Lambda$ ,  $|m_2 - m_{2c}| \leq \Lambda$ .

We introduce the  $Z$  variables

$$Z_{[i]}^{[J]} := (1 - \mathbb{E}_{[i]}) \left( G_{[ii]}^{[J]} \right)^{-1}.$$

By the identity (3.6) we have

$$G_{[ii]}^{-1} = \mathbb{E}_{[i]} G_{[ii]}^{-1} + Z_{[i]} = \begin{pmatrix} -w - w |d_i|^2 m_2^{[i]} & -w^{1/2} z \\ -w^{1/2} \bar{z} & -w - w m_1^{[i]} \end{pmatrix} + Z_{[i]}, \quad (4.9)$$

where

$$Z_{[i]} = w \begin{pmatrix} |d_i|^2 m_2^{[i]} - |d_i|^2 (XG^{[i]}X^\dagger)_{ii} & w^{-1/2} d_i X_{i\bar{i}} - (DXG^{[i]}DX)_{i\bar{i}} \\ w^{-1/2} \bar{d}_i X_{ii}^\dagger - (X^\dagger D^\dagger G^{[i]}X^\dagger D^\dagger)_{\bar{i}\bar{i}} & m_1^{[i]} - (X^\dagger D^\dagger G^{[i]}DX)_{\bar{i}\bar{i}} \end{pmatrix}. \quad (4.10)$$

**Lemma 4.3.** *For  $J \subseteq \mathcal{I}_1$ , the following crude bound on the difference between  $m_a$  and  $m_a^{[J]}$  ( $a = 1, 2$ ) holds:*

$$|m_a - m_a^{[J]}| \leq \frac{C|J|}{N\eta}, \quad a = 1, 2, \quad (4.11)$$

where  $C = C(\tau)$  is a constant depending only on  $\tau$ .

*Proof.* For  $i \in \mathcal{I}_1$ , we have

$$|m_1 - m_1^{(i)}| = \frac{1}{N} \left| \sum_{k \in \mathcal{I}_1} |d_k|^2 \frac{G_{ki} G_{ik}}{G_{ii}} \right| \leq \frac{\tau^{-1}}{N|G_{ii}|} \sum_{k \in \mathcal{I}_1} |G_{ik}|^2 = \frac{\tau^{-1}}{N\eta} \frac{\operatorname{Im} G_{ii}}{|G_{ii}|} \leq \frac{\tau^{-1}}{N\eta} \quad (4.12)$$

where in the first step we use (3.5), in the second and third steps the equality (3.15). Similarly, using (3.5) and (3.16) we get

$$|m_1^{(i)} - m_1^{(i\bar{i})}| = \frac{1}{N} \left| \sum_{k \in \mathcal{I}_1} |d_k|^2 \frac{G_{k\bar{i}}^{(i)} G_{ik}^{(i)}}{G_{i\bar{i}}^{(i)}} \right| \leq \frac{\tau^{-1}}{N|G_{i\bar{i}}^{(i)}|} \left( \frac{G_{i\bar{i}}^{(i)}}{|w|} + \frac{\bar{w}}{|w|} \frac{\operatorname{Im} G_{i\bar{i}}^{(i)}}{\eta} \right) \leq \frac{2\tau^{-1}}{N\eta}.$$

By induction on the indices in  $[J]$ , we can prove (4.12). The proof for  $m_2$  is similar.  $\square$

**Lemma 4.4.** *Suppose  $|z|^2 \leq 1 - \tau$ . For  $i \in \mathcal{I}_1$ , we have*

$$|(Z_{[i]})_{11}| < |w| \sqrt{\frac{\operatorname{Im} m_2^{[i]}}{N\eta}}, \quad |(Z_{[i]})_{22}| < |w| \sqrt{\frac{\operatorname{Im} m_1^{[i]}}{N\eta}}, \quad (4.13)$$

$$|(Z_{[i]})_{st}| < |w| \left( \frac{|w|^{-1/2}}{\sqrt{N}} + \sqrt{\frac{|m_1^{[i]}|}{N|w|}} + \sqrt{\frac{\text{Im } m_1^{[i]}}{N\eta}} \right) \text{ for } s \neq t \in \{1, 2\}, \quad (4.14)$$

uniformly in  $w \in \mathbf{D} \cup \mathbf{D}_L$ . In particular, these imply that

$$Z_{[i]} < |w|\Psi_\theta, \quad (4.15)$$

uniformly in  $w \in \mathbf{D}$ , and

$$Z_{[i]} < |w|(N\eta)^{-1/2}, \quad (4.16)$$

uniformly in  $w \in \mathbf{D}_L$ .

*Proof.* Apply the large deviation Lemma 3.6 to  $Z_{[i]}$  in (4.10), we get that

$$\begin{aligned} \left| \frac{(Z_{[i]})_{11}}{w} \right| &< \frac{1}{N} \left[ \left( \sum_{\mu} |G_{\mu\mu}^{[i]}|^2 \right)^{1/2} + \left( \sum_{\mu \neq \nu} |G_{\mu\nu}^{[i]}|^2 \right)^{1/2} \right] \leq \frac{C}{N} \left( \sum_{\mu, \nu} |G_{\mu\nu}^{[i]}|^2 \right)^{1/2} \\ &= \frac{C}{N} \left( \sum_{\mu} \frac{\text{Im } G_{\mu\mu}^{[i]}}{\eta} \right)^{1/2} = C \sqrt{\frac{\text{Im } m_2^{[i]}}{N\eta}}. \end{aligned}$$

where in the third step we use the equality (3.14). Similarly we can prove the bound for  $(Z_{[i]})_{22}$  using Lemma 3.6 and (3.15). Now we consider  $(Z_{[i]})_{12}$ . First, we have  $X_{i\bar{i}} < N^{-1/2}$  by (2.3). For the other part, we use Lemma 3.6 and (3.17) to get that

$$\begin{aligned} \left| (DXG^{[i]}DX)_{i\bar{i}} \right| &< \frac{1}{N} \left( \sum_{j, \mu} |d_j|^2 |G_{\mu j}^{[i]}|^2 \right)^{1/2} = \frac{1}{N} \left[ \sum_j |d_j|^2 \left( |w|^{-1} G_{jj}^{[i]} + \frac{\bar{w}}{|w|} \frac{\text{Im } G_{jj}^{[i]}}{\eta} \right) \right]^{1/2} \\ &\leq \left[ \frac{|m_1^{[i]}|}{N|w|} + \frac{\text{Im } m_1^{[i]}}{N\eta} \right]^{1/2} \leq C \left( \sqrt{\frac{|m_1^{[i]}|}{N|w|}} + \sqrt{\frac{\text{Im } m_1^{[i]}}{N\eta}} \right). \end{aligned} \quad (4.17)$$

Similarly we can prove the estimate for  $(Z_{[i]})_{21}$ .

Now we prove (4.15). By the definitions (4.7) and using (4.11), we get that

$$|(Z_{[i]})_{11}| < |w| \sqrt{\frac{\text{Im } m_2^{[i]}}{N\eta}} = |w| \sqrt{\frac{\text{Im } m_{2c} + \text{Im} (m_2^{[i]} - m_2) + \text{Im} (m_2 - m_{2c})}{N\eta}} \leq C|w|\Psi_\theta. \quad (4.18)$$

We can estimate  $(Z_{[i]})_{22}$  and the third term in (4.14) in a similar way. For the Cases 1-4 in Lemma 3.7, we have  $|m_{1c}| \sim 1$  for  $|w| \sim 1$ ,  $\text{Im } m_{1c} \sim |w|^{-1/2} \sim |m_{1c}|$  for  $|w| \rightarrow 0$ , and  $\eta \leq C\text{Im } m_{1c}$ . Thus

$$\sqrt{\frac{|m_{1c}|}{N|w|}} \leq \frac{C}{\sqrt{N}} \leq C\Psi_\theta \text{ for } |w| \sim 1, \quad \sqrt{\frac{|m_{1c}|}{N|w|}} \leq C\sqrt{\frac{\text{Im } m_{1c}}{N\eta}} \leq C\Psi_\theta \text{ for } |w| \rightarrow 0.$$

Then for the second term in (4.14), we have that

$$\sqrt{\frac{|m_1^{[i]}|}{N|w|}} \leq C \left( \frac{1}{N\eta} + \sqrt{\frac{\theta}{N\eta}} + \sqrt{\frac{|m_{1c}|}{N|w|}} \right) \leq C\Psi_\theta.$$

This concludes (4.15). Finally, the estimate (4.16) follows directly from (4.13), (4.14) and (3.13).  $\square$

**Lemma 4.5.** Suppose  $|z|^2 \leq 1 - \tau$ . Define the  $w$ -dependent event  $\Xi(w) := \{\theta \leq |w|^{-1/2}(\log N)^{-1}\}$ . Then we have that for  $w \in \mathbf{D}$ ,

$$\mathbf{1}(\Xi)m_2 = \mathbf{1}(\Xi) \left[ \frac{1 + m_1}{-w(1 + m_1)^2 + |z|^2} + O_{<}(\Psi_\theta) \right], \quad \mathbf{1}(\Xi)\Upsilon(w, m_1) < \mathbf{1}(\Xi)\Psi_\theta, \quad (4.19)$$

where  $\Upsilon$  is defined in (3.35). For  $w \in \mathbf{D}_L$ , we have

$$m_2 = \frac{1 + m_1}{-w(1 + m_1)^2 + |z|^2} + O_{<} \left( \eta^{-1}(N\eta)^{-1/2} \right), \quad \Upsilon(w, m_1) < \eta^{-1}(N\eta)^{-1/2}. \quad (4.20)$$

*Proof.* Using (4.9), we get

$$G_{[ii]}^{-1} = \pi_{[i]}^{-1} + \epsilon_{[i]}, \quad (4.21)$$

where  $\pi_{[i]}$  is defined in (2.36) and

$$\epsilon_{[i]} = w \begin{pmatrix} |d_i|^2 (m_2 - m_2^{[i]}) & 0 \\ 0 & m_1 - m_1^{[i]} \end{pmatrix} + Z_{[i]}.$$

By (4.11) and (4.15), we get that  $\epsilon_{[i]} < |w|\Psi_\theta$ . Let  $B_i = \pi_{[i]}^{-1} - \pi_{[i]c}^{-1}$ , where  $\pi_{[i]c}$  is defined in (2.32). By (3.31) and the definition of  $\Xi$ , we have  $\mathbf{1}(\Xi)\|B_i\pi_{[i]c}\| \leq C(\log N)^{-1}$ . Thus we have the expansion

$$\mathbf{1}(\Xi)\pi_{[i]} = \mathbf{1}(\Xi)(\pi_{[i]c}^{-1} + B_i)^{-1} = \mathbf{1}(\Xi)\pi_{[i]c} (1 - B_i\pi_{[i]c} + (B_i\pi_{[i]c})^2 + \dots) = \mathbf{1}(\Xi)(\pi_{[i]c} + \epsilon_a), \quad (4.22)$$

where  $\epsilon_a$  can be estimated as  $\mathbf{1}(\Xi)\|\epsilon_a\| \leq \mathbf{1}(\Xi)C|w|^{-1/2}(\log N)^{-1}$ . This shows that  $\mathbf{1}(\Xi)\|\pi_{[i]}\| = \mathbf{1}(\Xi)O(|w|^{-1/2})$ , and so  $\mathbf{1}(\Xi)\|\epsilon_{[i]}\pi_{[i]}\| < \mathbf{1}(\Xi)|w|^{1/2}\Psi_\theta \leq \mathbf{1}(\Xi)CN^{-\zeta/2}$  by the definition of  $\mathbf{D}$  in (2.39). Again we do the expansion for (4.21),

$$\mathbf{1}(\Xi)G_{[ii]} = \mathbf{1}(\Xi) \left( \pi_{[i]}^{-1} + \epsilon_{[i]} \right)^{-1} = \mathbf{1}(\Xi)\pi_{[i]} \left( 1 + \sum_{l=1}^{\infty} (-\epsilon_{[i]}\pi_{[i]})^l \right) = \mathbf{1}(\Xi) (\pi_{[i]} + \epsilon_b), \quad (4.23)$$

where  $\mathbf{1}(\Xi)\|\epsilon_b\| < \mathbf{1}(\Xi)\Psi_\theta$ . Now the 11 entry of (4.23) gives that

$$\mathbf{1}(\Xi)G_{ii} = \mathbf{1}(\Xi) \frac{-1 - m_1}{w(1 + |d_i|^2 m_2)(1 + m_1) - |z|^2} + \mathbf{1}(\Xi)O_{<}(\Psi_\theta), \quad (4.24)$$

from which we get that

$$\mathbf{1}(\Xi)G_{ii} \left[ -w(1 + |d_i|^2 m_2) + \frac{|z|^2}{1 + m_1} \right] = \mathbf{1}(\Xi) \left[ 1 + O_{<}(|w|^{1/2}\Psi_\theta) \right]. \quad (4.25)$$

Here we use that

$$\mathbf{1}(\Xi) \left[ -w(1 + |d_i|^2 m_2) + \frac{|z|^2}{1 + m_1} \right] = O(|w|^{1/2}),$$

which follows from Proposition 2.15 and the definition of  $\Xi$ . Summing (4.25) over  $i$ ,

$$\mathbf{1}(\Xi) \left[ -w(m_2 + m_1 m_2) + \frac{|z|^2 m_2}{1 + m_1} \right] = \mathbf{1}(\Xi) \left[ 1 + O_{<}(|w|^{1/2}\Psi_\theta) \right],$$

which gives

$$\mathbf{1}(\Xi)m_2 = \mathbf{1}(\Xi)\frac{1+m_1}{-w(1+m_1)^2+|z|^2} + \mathbf{1}(\Xi)O_{<}(\Psi_\theta). \quad (4.26)$$

Now plug (4.26) into (4.24), multiply with  $|d_i|^2$  and sum over  $i$ , we get

$$\mathbf{1}(\Xi)m_1 = \mathbf{1}(\Xi)\left[\frac{1}{N}\sum_{i=1}^n l_i s_i \frac{-1-m_1}{w\left(1+s_i\frac{1+m_1}{-w(1+m_1)^2+|z|^2}\right)(1+m_1)-|z|^2} + O_{<}(\Psi_\theta)\right], \quad (4.27)$$

where we use (3.29) and  $\mathbf{1}(\Xi)(1+m_1) = \mathbf{1}(\Xi)O(|w|^{-1/2})$ . This concludes the proof.

Similarly, when  $w \in \mathbf{D}_L$ , it is easy to prove (4.20) using the estimates (4.16) and (3.13). Note that  $|m_{1,2}| = O(\eta^{-1})$  by (3.13), which implies immediately the bounds  $\|\pi_{[i]}\| = O(\eta^{-1})$  and  $\|(\pi_{[i]})^{-1}\| = O(\eta)$ . Hence without introducing the event  $\Xi$ , we can obtain directly

$$G_{[ii]} = \pi_{[i]} + O_{<}(\eta^{-1}(N\eta)^{-1/2}). \quad (4.28)$$

The rest of the proof is essentially the same.  $\square$

Notice that applying Lemma 3.10 to (4.20), we obtain  $|m_{1,2} - m_{1,2c}| < \eta^{-1}(N\eta)^{-1/2}$ . Plugging it into (4.28), we get immediately (4.4) for  $w \in \mathbf{D}_L$ . This proves the entrywise law on  $\mathbf{D}_L$ , since  $\eta^{-1}N^{-1/2} \leq C\Psi$  by the definition (2.45) and the estimate (3.28).

## 4.2 The large $\eta$ case

It remains prove Proposition 4.1 on domain  $\mathbf{D}$ . We would like to fix  $E$  and then apply a continuity argument in  $\eta$  by first showing that the rough bound  $\Lambda \leq |w|^{-1/2}(\log N)^{-1}$  in Lemma 4.5 holds for large  $\eta$ . To start the argument, we first need to establish the estimates on  $G$  when  $\eta \sim 1$ . The next lemma is a trivial consequence of (3.13).

**Lemma 4.6.** *For any  $w \in \mathbf{D}$  and  $\eta \geq c$  for fixed  $c > 0$ , we have the bound*

$$\max_{s,t} |G_{st}(w)| \leq C \quad (4.29)$$

for some  $C > 0$ . This estimate also holds if we replace  $G$  with  $G^{(J)}$  for  $J \subset \mathcal{I}$ .

**Lemma 4.7.** *Fix  $c > 0$  and  $|z|^2 \leq 1 - \tau$ . We have the following estimate*

$$\max_{w \in \mathbf{D}, \eta \geq c} \Lambda(w) < N^{-1/2}. \quad (4.30)$$

*Proof.* By the previous lemma, we have  $|m_{1,2}^{[i]}| = O(1)$ . So by Lemma 4.4,  $\|Z_{[i]}\| < N^{-1/2}$  uniformly in  $\eta \geq c$ . Then as in (4.21),

$$G_{[ii]} = \left(\pi_{[i]}^{-1} + \epsilon_{[i]}\right)^{-1}, \quad (4.31)$$

where  $\|\pi_{[i]}^{-1}\| = O(1)$  and  $\|\epsilon_{[i]}\| < N^{-1/2}$ . Notice since  $G_{[ii]} = O(1)$ , we have the estimate

$$\pi_i = \left(G_{[ii]}^{-1} - \epsilon_{[i]}\right)^{-1} = G_{[ii]}(1 - \epsilon_{[i]}G_{[ii]})^{-1} = O_{<}(1).$$

Then we can expand (4.31) to get that

$$G_{[ii]} = \pi_i + O_{<}(N^{-1/2}). \quad (4.32)$$

The 11 and 22 entries of (4.32) leads to the equations

$$m_1 = \frac{1}{N} \sum_{i=1}^N |d_i|^2 \left[ -w(1 + |d_i|^2 m_2) + \frac{|z|^2}{1 + m_1} \right]^{-1} + O_{<}(N^{-1/2}), \quad (4.33)$$

$$m_2 = \frac{1}{N} \sum_{i=1}^N \left[ -w(1 + m_1) + \frac{|z|^2}{1 + |d_i|^2 m_2} \right]^{-1} + O_{<}(N^{-1/2}). \quad (4.34)$$

Our goal is to prove that  $\text{Im } m_{1,2} \geq C(\log N)^{-1}$  with high probability for some  $C > 0$ .

Using the spectral decomposition (3.11), we note that for  $l > 1$ ,

$$\begin{aligned} \frac{1}{N} \sum_{|\lambda_k - E| \geq l\eta} \frac{|E - \lambda_k|}{(\lambda_k - E)^2 + \eta^2} &\leq \frac{1}{l\eta}, \\ \frac{1}{N} \sum_{|\lambda_k - E| \leq l\eta} \frac{|E - \lambda_k|}{(\lambda_k - E)^2 + \eta^2} &\leq \frac{1}{N} \sum_{|\lambda_k - E| \leq l\eta} \frac{l\eta}{(\lambda_k - E)^2 + \eta^2} \leq l \text{Im } m_2. \end{aligned}$$

Summing up these two inequalities and optimizing  $l$ , we get

$$|\text{Re } m_2| \leq 2\sqrt{\frac{\text{Im } m_2}{\eta}}. \quad (4.35)$$

Assume that  $\text{Im } m_2 \leq C(\log N)^{-1}$ , then by (4.8) we also have  $\text{Im } m_1 \leq C\tau^{-1}(\log N)^{-1}$ . From (4.35), we get  $|m_2| \leq C(\log N)^{-1/2}$ . Together with  $\text{Im } w = \eta \geq c$  and  $\text{Im}[|z|^2/(1 + m_1)] < 0$ , (4.33) gives

$$|m_1| \leq \frac{1}{N} \sum_i |d_i|^2 \left| \text{Im} \left[ -w(1 + |d_i|^2 m_2) + \frac{|z|^2}{1 + m_1} \right] \right|^{-1} + o(1) \leq C \quad (4.36)$$

with high probability. Using the above estimate and  $|m_2| \leq C(\log N)^{-1/2}$  we get

$$\left| -w(1 + m_1) + \frac{|z|^2}{1 + |d_i|^2 m_2} \right| \leq C \text{ with high probability.}$$

On the other hand

$$\text{Im} \left[ -w(1 + m_1) + \frac{|z|^2}{1 + |d_i|^2 m_2} \right] \leq -\text{Im } w = -\eta, \quad (4.37)$$

where we use  $\text{Im}[|z|^2/(1 + |d_i|^2 m_2)] < 0$  and

$$\text{Im}(wm_1) = \text{Im} \left[ \frac{1}{N} \sum_{k=1}^N |d_i|^2 |\xi_k(i)|^2 \left( -1 + \frac{\lambda_k}{\lambda_k - w} \right) \right] \geq 0.$$

Hence (4.34) implies  $\text{Im } m_2 \geq c$  with high probability for some  $c > 0$ . This contradicts  $\text{Im } m_2 \leq C(\log N)^{-1}$ . Thus  $\text{Im } m_2 \geq C(\log N)^{-1}$  with high probability for some  $C > 0$ , which also implies  $\text{Im } m_1 \geq C(\log N)^{-1}$  by (4.8).

Now we can proceed as in Lemma 4.5 and get that

$$m_2 = \frac{1 + m_1}{-w(1 + m_1)^2 + |z|^2} + O_{<}(N^{-1/2}), \quad \Upsilon(w, m_1) < N^{1/2}. \quad (4.38)$$

We omit the details. Applying Lemma 3.10 to (4.38), we conclude  $|m_{1,2} - m_{1,2c}| < N^{-1/2}$  uniformly in  $\eta \geq c$ . By (4.32), we get  $\|(G - \Pi)_{[ii]}\| < N^{-1/2}$  uniformly in  $\eta \geq c$  and  $i \in \mathcal{I}_1$ . Finally using (3.8) and Lemmas 3.5-3.6, we can prove the off-diagonal estimate (see (4.51)).  $\square$

### 4.3 Proof of the weak entrywise local law

In this subsection, we finish the proof of Proposition 4.1 on domain  $\mathbf{D}$ . We shall fix the real part  $E$  of  $w = E + i\eta$  and decrease the imaginary part  $\eta$ . Recall Lemma 4.5 is based on the condition  $\Lambda \leq |w|^{-1/2}(\log N)^{-1}$  (i.e. event  $\Xi$ ). So far this is only established for large  $\eta$  in (4.30). We want to show this condition for small  $\eta$  also by using a continuity argument.

It is convenient to introduce the random function

$$v(w) = \max_{w' \in L(w)} \theta(w') |w'|^{1/2} \left( \frac{N \operatorname{Im} w'}{|w'|^{1/2}} \right)^{1/4},$$

where  $L(w)$  is defined in (3.36). Fix a regular domain  $\mathbf{S}$ , an  $\epsilon < \zeta/4$  and a large  $D > 0$ . Our goal is to prove that with high probability there is a gap in the range of  $v$ , i.e.

$$\mathbb{P} \left( v(w) \leq N^\epsilon, v(w) > N^{3\epsilon/4} \right) \leq N^{-D+21} \quad (4.39)$$

for all  $w \in \mathbf{S}$  and large enough  $N \geq N(\epsilon, D)$ .

Suppose  $v(w) \leq N^\epsilon$ , then it is easy to verify

$$\theta(w') \leq C |w'|^{-1/2} (\log N)^{-1} \quad (4.40)$$

for all  $w' \in L(w)$ . Hence  $\{v(w) \leq N^\epsilon\} \subset \Xi(w')$  for all  $w' \in \mathbf{S} \cap L(w)$ . Then by (4.19), we have that for all  $w' \in \mathbf{S} \cap L(w)$ , there exists an  $N_0 \equiv N_0(\epsilon, D)$  such that

$$P \left( v(w) \leq N^\epsilon, \Upsilon(w') > \frac{N^\epsilon}{|w'|^{1/2}} \sqrt{\frac{|w'|^{1/2}}{N \operatorname{Im} w'}} \right) \leq N^{-D}, \quad (4.41)$$

for all  $N > N_0$ . Taking the union bound we get

$$P \left( v(w) \leq N^\epsilon, \max_{w' \in L(w)} \Upsilon(w') \sqrt{\frac{N \operatorname{Im} w'}{|w'|^{-1/2}}} > N^\epsilon \right) \leq N^{-D+10}. \quad (4.42)$$

Now consider the event

$$\Xi_1 := \left\{ v(w) \leq N^\epsilon, \max_{w' \in L(w)} \Upsilon(w') \sqrt{\frac{N \operatorname{Im} w'}{|w'|^{-1/2}}} \leq N^\epsilon \right\}. \quad (4.43)$$

Then  $1(\Xi_1) \Upsilon(w') \leq \delta(w')$  for all  $w' \in L(w)$  with  $\delta(w') = \frac{N^\epsilon}{|w'|^{1/2}} \sqrt{\frac{|w'|^{1/2}}{N \operatorname{Im} w'}}$ . We now apply Lemma 3.10. If  $\kappa \ll 1$  (recall (3.21)), then  $|w| \sim 1$  and we have

$$1(\Xi_1) |m_1(w') - m_{1c}(w')| \leq C \sqrt{\delta(w')} \leq C N^{\epsilon/2} \left( \frac{1}{N \operatorname{Im} w'} \right)^{1/4}$$

for all  $w' \in L(w)$ ; if  $\kappa \geq c > 0$  for some constant  $c > 0$ , then

$$1(\Xi_1)|m_1(w') - m_{1c}(w')| \leq C\delta(w') \leq C \frac{N^\epsilon}{|w'|^{1/2}} \left( \frac{|w'|^{1/2}}{N \operatorname{Im} w'} \right)^{1/2}$$

for all  $w' \in L(w)$ . Combining these two cases we get

$$1(\Xi_1)|m_1(w') - m_{1c}(w')| \leq C \frac{N^{\epsilon/2}}{|w'|^{1/2}} \left( \frac{|w'|^{1/2}}{N \operatorname{Im} w'} \right)^{1/4} \quad (4.44)$$

for all  $w' \in L(w)$ . By (4.19), we have

$$1(\Xi_1)|m_2(w') - m_{2c}(w')| < 1(\Xi_1)|m_1(w') - m_{1c}(w')| + 1(\Xi_1)\Psi_\theta < \frac{N^{\epsilon/2}}{|w'|^{1/2}} \left( \frac{|w'|^{1/2}}{N \operatorname{Im} w'} \right)^{1/4},$$

for all  $w' \in \mathbf{S} \cap L(w)$ . Combining this bound with (4.44), we see there is  $N_1 \equiv N_1(\epsilon, D)$  such that

$$\mathbb{P} \left( v(w) \leq N^\epsilon, \max_{w' \in L(w)} \Upsilon(w') \sqrt{\frac{N \operatorname{Im} w'}{|w'|^{-1/2}}} \leq N^\epsilon, \max_{w' \in L(w)} \theta(w') |w'|^{1/2} \left( \frac{N \operatorname{Im} w'}{|w'|^{1/2}} \right)^{1/4} > N^{3\epsilon/4} \right) \leq N^{-D} \quad (4.45)$$

for  $N \geq \max\{N_0, N_1\}$ . Adding (4.42) and (4.45), we get

$$\mathbb{P} \left( v(w) \leq N^\epsilon, \max_{w' \in L(w)} \theta(w') |w'|^{1/2} \left( \frac{N \operatorname{Im} w'}{|w'|^{1/2}} \right)^{1/4} > N^{3\epsilon/4} \right) \leq N^{-D+11}.$$

Taking the union bound over  $L(w)$  we get (4.39) for all  $N \geq \max\{N_0, N_1\}$ .

Now we conclude the proof of Proposition 4.1 by combining (4.39) with the large  $\eta$  estimate (4.30). We choose a lattice  $\Delta \subset \mathbf{S}$  such that  $|\Delta| \leq N^{20}$  and for any  $w \in \mathbf{S}$  there is a  $w' \in \Delta$  with  $|w' - w| \leq N^{-9}$ . Taking the union bound we get

$$\mathbb{P} \left( \exists w \in \Delta : v(w) \in (N^{3\epsilon/4}, N^\epsilon] \right) \leq N^{-D+41}. \quad (4.46)$$

Since  $v$  has Lipschitz constant bounded by, say,  $N^6$ , then we have

$$\mathbb{P} \left( \exists w \in \mathbf{S} : v(w) \in (2N^{3\epsilon/4}, N^\epsilon/2] \right) \leq N^{-D+41}. \quad (4.47)$$

Combining with (4.30), we see that there exists  $N_2 \equiv N_2(\epsilon, D)$  such that for  $N > N_2$ ,

$$\mathbb{P} \left( \forall w \in \mathbf{S} : v(w) \leq 2N^{3\epsilon/4} \right) \geq 1 - 2N^{-D+41}.$$

Since  $\epsilon$  and  $D$  are arbitrary, the above inequality shows that  $v(w) < 1$  uniformly in  $w \in \mathbf{S}$ , or

$$\theta(w) < \frac{1}{|w|^{1/2}} \left( \frac{|w|^{1/2}}{N\eta} \right)^{1/4}. \quad (4.48)$$

In particular we see that for all  $w \in \mathbf{S}$ , the event  $\Xi$  holds with high-probability.

Now using (4.23) and (4.48), we get

$$\|G_{[ii]} - \pi_{[i]c}\| \leq \|G_{[ii]} - \pi_{[i]}\| + \|\pi_{[i]} - \pi_{[i]c}\| < \Psi_\theta + \theta < \frac{1}{|w|^{1/2}} \left( \frac{|w|^{1/2}}{N\eta} \right)^{1/4}. \quad (4.49)$$

To conclude Proposition 4.1, it remains to prove the estimate for the off-diagonal entries. By (4.11), it is not hard to see that

$$\|G_{[ii]}^{[J]} - \pi_{[i]c}\| < \frac{1}{|w|^{1/2}} \left( \frac{|w|^{1/2}}{N\eta} \right)^{1/4} \quad (4.50)$$

for any  $|J| \leq l$  with  $l \in \mathbb{N}$  fixed. Thus we have  $G_{[ii]}^{[J]} = O(|w|^{-1/2})$  and  $(G_{[ii]}^{[J]})^{-1} = O(|w|^{1/2})$  with high probability. Let  $i \neq j \in \mathcal{I}_1$ , using (3.8) and the above diagonal estimates, we get that

$$\|G_{[ij]}\| < |w|^{-1} \frac{|w|^{1/2}}{\sqrt{N}} + |w|^{-1} \left\| \sum_{k, l \notin \{i, j\}} H_{[ik]} G_{[kl]}^{[ij]} H_{[lj]} \right\| < \Psi_\theta < \frac{1}{|w|^{1/2}} \left( \frac{|w|^{1/2}}{N\eta} \right)^{1/4}, \quad (4.51)$$

where, as in the proof of Lemma 4.4, we use Lemmas 3.5 and 3.6 to obtain that

$$|w|^{-1} \left\| \sum_{k, l \notin \{i, j\}} H_{[ik]} G_{[kl]}^{[ij]} H_{[lj]} \right\| = \left\| \begin{pmatrix} \sum_{k, l \notin \{i, j\}} X_{i\bar{k}} G_{kl}^{[ij]} X_{l\bar{j}}^\dagger & \sum_{k, l \notin \{i, j\}} X_{i\bar{k}} G_{kl}^{[ij]} X_{l\bar{j}} \\ \sum_{k, l \notin \{i, j\}} X_{i\bar{k}}^\dagger G_{kl}^{[ij]} X_{l\bar{j}}^\dagger & \sum_{k, l \notin \{i, j\}} X_{i\bar{k}}^\dagger G_{kl}^{[ij]} X_{l\bar{j}} \end{pmatrix} \right\| < \Psi_\theta. \quad (4.52)$$

#### 4.4 Proof of the strong entrywise local law

In this section, we finish the proof of the (strong) entrywise local law in Theorem 2.18 on domain  $\mathbf{D}$  and under the condition  $|w|^{1/2} + |z|^2 \geq c$ . In Lemma 4.5, we have proved an error estimate of the self-consistent equations of  $m_{1,2}$  linearly in  $\Psi_\theta$ . The core part of the proof is to improve this estimate to quadratic in  $\Psi_\theta$ . For the sequence of random variables  $Z_{[i]}$ , we define the averaged quantities

$$[Z] = \frac{1}{N} \sum_{i=1}^N \pi_{[i]} Z_{[i]} \pi_{[i]}, \quad \langle Z \rangle = \frac{1}{N} \sum_{i=1}^N |d_i|^2 \pi_{[i]} Z_{[i]} \pi_{[i]}.$$

The following Lemma is an improvement of Lemma 4.5.

**Lemma 4.8.** *Fix  $|z|^2 \leq 1 - \tau$ . Then for  $w \in \mathbf{D}$ ,*

$$m_2 = \frac{1 + m_1}{-w(1 + m_1)^2 + |z|^2} + O_{<}(|w|^{1/2} \Psi_\theta^2 + \|[Z]\| + \|\langle Z \rangle\|), \quad (4.53)$$

and

$$\Upsilon(w, m_1) < |w|^{1/2} \Psi_\theta^2 + \|[Z]\| + \|\langle Z \rangle\|. \quad (4.54)$$

For  $w \in \mathbf{D}_L$ ,

$$m_2 = \frac{1 + m_1}{-w(1 + m_1)^2 + |z|^2} + O_{<}((N\eta)^{-1} + \|[Z]\| + \|\langle Z \rangle\|), \quad (4.55)$$

and

$$\Upsilon(w, m_1) < (N\eta)^{-1} + \|[Z]\| + \|\langle Z \rangle\|. \quad (4.56)$$



*Proof.* The proof is almost the same as the one in Lemma 4.5, we only lay out the difference. We first consider the case  $w \in \mathbf{D}$ . By Proposition 4.1, the event  $\Xi$  holds with high probability. Hence without loss of generality, we may assume  $\Xi$  holds throughout the proof. Using (3.9), we get

$$\frac{1}{N} \sum_{k \in \mathcal{I}_1} \begin{pmatrix} |d_k|^2 & 0 \\ 0 & 1 \end{pmatrix} (G_{[kk]} - G_{[kk]}^{[i]}) = \begin{pmatrix} |d_i|^2 & 0 \\ 0 & 1 \end{pmatrix} \frac{G_{[ii]}}{N} + \frac{1}{N} \sum_{k \neq i} \begin{pmatrix} |d_k|^2 & 0 \\ 0 & 1 \end{pmatrix} G_{[ki]} G_{[ii]}^{-1} G_{[ik]}. \quad (4.57)$$

By Proposition 4.1, (3.31) and (4.51), we have

$$\|G_{[ki]} G_{[ii]}^{-1} G_{[ik]}\| < |w|^{1/2} \Psi_\theta^2.$$

By Lemma 3.7, it is easy to verify that  $\|G_{[ii]}/N\| \leq C|w|^{1/2} \Psi_\theta^2$ . Plug it into (4.57), we get

$$|m_{1,2}^{[i]} - m_{1,2}| < |w|^{1/2} \Psi_\theta^2. \quad (4.58)$$

Using (4.15) and (4.58), the error  $\epsilon_b$  in (4.23) is

$$\epsilon_b = O_{<}(|w|^{1/2} \Psi_\theta^2) - \pi_{[i]} Z_{[i]} \pi_{[i]} \left[ 1 + O_{<}(|w|^{1/2} \Psi_\theta) \right] = O_{<}(|w|^{1/2} \Psi_\theta^2) - \pi_{[i]} Z_{[i]} \pi_{[i]}.$$

Then following the arguments in Lemma 4.5, we can obtain the desired result on  $\Xi$ . For  $w \in \mathbf{D}_L$ , the proof is similar by using (4.4).  $\square$

In the following lemma we prove stronger bounds on  $[Z]$  and  $\langle Z \rangle$  by keeping track of the cancellation effects due to the average over the index  $i$ . The proof is given in Appendix B.

**Lemma 4.9.** (*Fluctuation averaging*) Fix  $|z|^2 \leq 1 - \tau$ . Suppose  $\Phi$  and  $\Phi_o$  are positive,  $N$ -dependent deterministic functions satisfying  $N^{-1/2} \leq \Phi, \Phi_o \leq N^{-c}$  for some constant  $c > 0$ . Suppose moreover that  $\Lambda < |w|^{-1/2} \Phi$  and  $\Lambda_o < |w|^{-1/2} \Phi_o$ . Then for  $w \in \mathbf{D}$ ,

$$\|[Z]\| + \|\langle Z \rangle\| < |w|^{-1/2} \Phi_o^2. \quad (4.59)$$

Now we finish the proof of the entrywise local law and averaged local law on the domain  $\mathbf{D}$ . By Proposition 4.1, we can take in Lemma 4.9

$$\Phi_o = |w|^{1/2} \sqrt{\frac{\text{Im}(m_{1c} + m_{2c}) + |w|^{-3/8} (N\eta)^{-1/4}}{N\eta}}, \quad \Phi = \left( \frac{|w|^{1/2}}{N\eta} \right)^{1/4},$$

with  $\Lambda_o < \Psi_\theta < |w|^{-1/2} \Phi_o$  and  $\Lambda < \theta < |w|^{-1/2} \Phi$ . Then (4.54) gives

$$\Upsilon(w, m_1) < \frac{|w|^{1/2} \text{Im}(m_{1c} + m_{2c}) + |w|^{1/4} (N\eta)^{-1/4}}{N\eta}.$$

Then using the stability Lemma 3.10,

$$|m_1 - m_{1c}| < \frac{|w|^{1/2} \text{Im}(m_{1c} + m_{2c})}{N\eta \sqrt{\kappa + \eta}} + \frac{|w|^{1/8}}{(N\eta)^{5/8}} < \frac{1}{N\eta} + \frac{|w|^{1/8}}{(N\eta)^{5/8}} < |w|^{-1/2} \left( \frac{|w|^{1/2}}{N\eta} \right)^{1/2+1/8}.$$

Here if  $\sqrt{\kappa + \eta} \geq (\log N)^{-1}$ , we use

$$\frac{|w|^{1/2} \text{Im}(m_{1c} + m_{2c})}{N\eta \sqrt{\kappa + \eta}} \leq \frac{C \log N}{N\eta} < \frac{1}{N\eta},$$

while if  $\sqrt{\kappa + \eta} \leq (\log N)^{-1}$ , we have  $\text{Im}(m_{1c} + m_{2c}) = O(\sqrt{\kappa + \eta})$ , which also gives that

$$\frac{|w|^{1/2} \text{Im}(m_{1c} + m_{2c})}{N\eta\sqrt{\kappa + \eta}} < \frac{1}{N\eta}.$$

We then use (4.53) to get that

$$\theta < |m_1 - m_{1c}| + \frac{|w|^{1/2} \text{Im}(m_{1c} + m_{2c}) + |w|^{1/4} (N\eta)^{-1/4}}{N\eta} < |w|^{-1/2} \left( \frac{|w|^{1/2}}{N\eta} \right)^{1/2+1/8}. \quad (4.60)$$

Repeating the previous steps with the new estimate (4.60), we get the bound

$$\theta < |w|^{-1/2} \left( \frac{|w|^{1/2}}{N\eta} \right)^{\sum_{k=1}^l 1/2^k + 1/2^{l+2}}$$

after  $l$  iterations. This implies the averaged local law  $\theta < (N\eta)^{-1}$  since  $l$  can be arbitrarily large. Finally as in (4.49) and (4.51), we have for  $i \neq j$

$$\|G_{[ii]} - \pi_{[i]c}\| + \|G_{[ij]}\| < \Psi_\theta + \theta < \sqrt{\frac{\text{Im}(m_{1c} + m_{2c})}{N\eta}} + \frac{1}{N\eta}.$$

This concludes the entrywise local law and averaged local law in Theorem 2.18 when  $|w|^{1/2} + |z|^2 \sim 1$ .

When  $w \in \mathbf{D}_L$ , we have proved the entrywise law (see the remark after (4.28)). Also we can prove a similar result as Lemma 4.9, which implies

$$m_2 = \frac{1 + m_1}{-w(1 + m_1)^2 + |z|^2} + O_{<}((N\eta)^{-1}), \quad \Upsilon(w, m_1) < (N\eta)^{-1}. \quad (4.61)$$

The averaged local law then follows from Lemma 3.10. We leave the details to the reader.

## 4.5 Proof of Theorem 2.18 when $|z|$ and $|w|$ are small

In the previous proof, we did not include the case where  $|w|^{1/2} + |z|^2 \leq \epsilon$  for some sufficiently small constant  $\epsilon > 0$ . The only reason is that Lemma 3.10 does not apply in this case. In this section, we deal with this problem.

The main idea of this subsection is to use a different set of self-consistent equations, which has the desired stability when  $|w|$  and  $|z|$  are small. Multiplying (4.24) with  $|d_i|^2$  and summing over  $i$ ,

$$1(\Xi)m_1 = 1(\Xi) \left[ \frac{1}{N} \sum_{i=1}^n l_i s_i \frac{-1 - m_1}{w(1 + s_i m_2)(1 + m_1) - |z|^2} + O_{<}(\Psi_\theta) \right]. \quad (4.62)$$

Recall that  $\Sigma := DD^\dagger = D^\dagger D$ . We introduce a new matrix

$$\tilde{H}(w) := \begin{pmatrix} -w\Sigma^{-1} & w^{1/2}(X - D^{-1}z) \\ w^{1/2}(X - D^{-1}z)^\dagger & -wI \end{pmatrix}, \quad (4.63)$$

and define  $\tilde{G} := \tilde{H}^{-1}$ . By Schur's complement formula, the upper left block of  $\tilde{G}$  is

$$\tilde{G}_L = [(X - D^{-1}z)(X - D^{-1}z)^\dagger - w\Sigma^{-1}]^{-1},$$

and the lower right block is equal to

$$\tilde{G}_R = [(X - D^{-1}z)^\dagger \Sigma (X - D^{-1}z) - w]^{-1} = [(DX - z)^\dagger (DX - z) - w]^{-1} = G_R.$$

Now we write  $m_{1,2}$  in another way as

$$m_1 = \frac{1}{N} \text{Tr} \left[ D^\dagger (Y Y^\dagger - w)^{-1} D \right] = \frac{1}{N} \text{Tr} \tilde{G}_L, \quad (4.64)$$

$$\begin{aligned} m_2 &= \frac{1}{N} \text{Tr} \tilde{G}_R = \frac{1}{N} \text{Tr} [(X - D^{-1}z)^\dagger \Sigma (X - D^{-1}z) - w]^{-1} \\ &= \frac{1}{N} \text{Tr} [(X - D^{-1}z)(X - D^{-1}z)^\dagger \Sigma - w]^{-1} = \frac{1}{N} \text{Tr} (\Sigma^{-1} \tilde{G}_L). \end{aligned} \quad (4.65)$$

We apply the arguments in the proof of Lemma 4.5 to  $\tilde{H}$ , and get that

$$\tilde{G}_{[ii]}^{-1} = \begin{pmatrix} -w|d_i|^{-2} - wm_2 & -w^{1/2}z\bar{d}_i^{-1} \\ -w^{1/2}\bar{z}\bar{d}_i^{-1} & -w - wm_1 \end{pmatrix} + O_{<}(|w|\Psi_\theta), \quad (4.66)$$

from which we get that

$$1(\Xi)\tilde{G}_{ii} = 1(\Xi) \left[ \frac{-1 - m_1}{w(|d_i|^{-2} + m_2)(1 + m_1) - |z|^2|d_i|^{-2}} + O_{<}(\Psi_\theta) \right].$$

Plugging this into (4.65), we get

$$1(\Xi)m_2 = 1(\Xi) \left[ \frac{1}{N} \sum_{i=1}^n \frac{l_i}{s_i} \frac{-1 - m_1}{w(s_i^{-1} + m_2)(1 + m_1) - |z|^2 s_i^{-1}} + O_{<}(\Psi_\theta) \right]. \quad (4.67)$$

We take the equations in (4.62) and (4.67) as our new self-consistent equations, namely,

$$1(\Xi)f_1(m_1, m_2) = 1(\Xi)O(\Psi_\theta), \quad 1(\Xi)f_2(m_1, m_2) = 1(\Xi)O(\Psi_\theta), \quad (4.68)$$

where

$$f_1(m_1, m_2) := m_1 + \frac{1}{N} \sum_i l_i s_i \frac{1 + m_1}{w(1 + s_i m_2)(1 + m_1) - |z|^2}, \quad (4.69)$$

$$f_2(m_1, m_2) := m_2 + \frac{1}{N} \sum_i l_i \frac{1 + m_1}{w(1 + s_i m_2)(1 + m_1) - |z|^2}. \quad (4.70)$$

According to the following lemma, this system of self-consistent equations are stable when  $|w|$  and  $|z|^2$  are small enough .

**Lemma 4.10.** *Suppose that  $N^{-2}|w|^{-1/2} \leq \delta(w) \leq (\log N)^{-1}|w|^{-1/2}$  for  $w \in \mathbf{D}$ . Suppose  $u_{1,2} : \mathbf{D} \rightarrow \mathbb{C}$  are Stieltjes transforms of positive integrable functions such that*

$$\max \{ |f_1(u_1, u_2)(w)|, |f_2(u_1, u_2)(w)| \} \leq \delta(w).$$

*Then there exists an  $\epsilon > 0$  such that if  $|w|^{1/2} + |z|^2 \leq \epsilon$ , we have*

$$|u_1(w) - m_{1c}(w)| + |u_2(w) - m_{2c}(w)| \leq C\delta, \quad (4.71)$$

*for some constant  $C > 0$  independent of  $w$ ,  $z$  and  $N$ .*

*Proof.* The proof depends on the estimate of the Jacobian at  $(m_{1c}, m_{2c})$ . By (3.26) and (A.35),

$$m_{1c} = \frac{i\sqrt{t_0} + O(|w|^{1/2} + |z|^2)}{\sqrt{w}}, \quad m_{2c} = \frac{it_0^{-1/2} + O(|w|^{1/2} + |z|^2)}{\sqrt{w}},$$

where  $t_0 = (N^{-1} \sum_{i=1}^n l_i/s_i)^{-1}$ . Then we can calculate that

$$\det \begin{pmatrix} \partial_1 f_1 & \partial_2 f_1 \\ \partial_1 f_2 & \partial_2 f_2 \end{pmatrix}_{u_{1,2}=m_{1,2c}} = \det \begin{pmatrix} 1 + O(|z|^2) & t_0 + O(|w|^{1/2} + |z|^2) \\ O(|z|^2) & 2 + O(|w|^{1/2} + |z|^2) \end{pmatrix} = 2 + O(|w|^{1/2} + |z|^2).$$

We can conclude the stability by expanding  $f_{1,2}(u_1, u_2)$  around  $(m_{1c}, m_{2c})$  and using a fixed point argument as in the proof of Lemma 3.10 in Section A.3.  $\square$

With this stability lemma, we can repeat all the arguments in the previous subsections to prove the entrywise local law and averaged local law when  $|w|^{1/2} + |z|^2 \leq \epsilon$ .

## 5 Anisotropic local law when $T$ is diagonal

In this section we prove the anisotropic local law in Theorem 2.18 when  $T$  is diagonal. The basic ideas of the proof follow from [4, section 5], and the core part of our proof is a novel way to perform the combinatorics. By the Definition 2.17 (ii) and the definition of matrix norm, it suffices to prove the following proposition for generalized entries of  $G$ .

**Proposition 5.1.** *Fix  $|z|^2 \leq 1 - \tau$  and suppose that the assumptions of Theorem 2.18 hold. Then for any regular domain  $\mathbf{S}$ ,*

$$|\langle \mathbf{u}, (G(w) - \Pi(w)) \mathbf{v} \rangle| < \Psi \quad (5.1)$$

*uniformly in  $w \in \mathbf{S}$  and any deterministic unit vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{C}^{\mathcal{I}}$ .*

It is equivalent to show that

$$\sum_{i,j \in \mathcal{I}_1} u_{[i]}^\dagger (G_{[ij]} - \Pi_{[ij]}) v_{[j]} < \Psi, \quad u_{[i]} := \begin{pmatrix} u_i \\ u_{\bar{i}} \end{pmatrix}, \quad v_{[j]} := \begin{pmatrix} v_j \\ v_{\bar{j}} \end{pmatrix}. \quad (5.2)$$

By the entrywise local law,

$$\left| \sum_{i,j} u_{[i]}^\dagger (G_{[ij]} - \Pi_{[ij]}) v_{[j]} \right| \leq \sum_i \|G_{[ii]} - \Pi_{[ii]}\| |u_{[i]}| |v_{[i]}| + \left| \sum_{i \neq j} u_{[i]}^\dagger G_{[ij]} v_{[j]} \right| < \Psi + \left| \sum_{i \neq j} u_{[i]}^\dagger G_{[ij]} v_{[j]} \right|.$$

Thus to show (5.2), it suffices to prove

$$\left| \sum_{i \neq j} u_{[i]}^\dagger G_{[ij]} v_{[j]} \right| < \Psi. \quad (5.3)$$

Notice from the entrywise law, we can only get

$$\left| \sum_{i \neq j} u_{[i]}^\dagger G_{[ij]} v_{[j]} \right| < \Psi \|\mathbf{u}\|_1 \|\mathbf{v}\|_1 \leq N\Psi,$$

using  $\|\mathbf{u}\|_1 \leq N^{1/2} \|\mathbf{u}\|_2$  and  $\|\mathbf{v}\|_1 \leq N^{1/2} \|\mathbf{v}\|_2$ . In particular, this estimate of the  $\ell^1$  norm is sharp when  $\mathbf{u}, \mathbf{v}$  are delocalized, i.e. their entries have size of order  $N^{-1/2}$ .

The estimate (5.3) follows from the Chebyshev's inequality if we can prove the following lemma.

**Lemma 5.2.** *Suppose the assumptions in Proposition 5.1 hold. For any even  $p \in 2\mathbb{N}$ , there exists a constant  $C_p$  which is independent of  $N$  such that*

$$\mathbb{E} \left| \sum_{i \neq j} u_{[i]}^\dagger G_{[ij]} v_{[j]} \right|^p \leq C_p \Psi^p.$$

The proof of Lemma 5.2 is based on the polynomialization method developed in [4, section 5]. Again we only give the proof for  $w \in \mathbf{D}$ . When  $w \in \mathbf{D}_L$ , the proof is almost the same.

### 5.1 Rescaling and partition of indices

For our purpose, it is convenient to define the rescaled matrix

$$R^{(J)} := w^{1/2} G^{(J)}, \quad (5.4)$$

for any  $J \subset \mathcal{I}$  and  $|J| \leq l$  for some fixed  $l$ . Consequently we define the control parameter  $\Phi$

$$\Phi = |w|^{1/2} \Psi. \quad (5.5)$$

By the entrywise law, for  $w \in \mathbf{D}$ ,

$$R_{[ii]}^{(J)} = O_{<}(1), \quad \left( R_{[ii]}^{(J)} \right)^{-1} = O_{<}(1), \quad R_{[ij]}^{(J)} = O_{<}(\Phi) \text{ for } i \neq j \quad (5.6)$$

under the above scaling. Now to prove Lemma 5.2, it is equivalent to prove

$$\mathbb{E} \left| \sum_{i \neq j} u_{[i]}^\dagger R_{[ij]} v_{[j]} \right|^p \leq C_p \Phi^p. \quad (5.7)$$

We expand the product in (5.7) as

$$\left| \sum_{i \neq j} u_{[i]}^\dagger R_{[ij]} v_{[j]} \right|^p = \sum_{i_k \neq j_k \in \mathcal{I}_1} \prod_{k=1}^{p/2} u_{[i_k]}^\dagger R_{[i_k j_k]} v_{[j_k]} \cdot \prod_{k=p/2+1}^p \overline{u_{[i_k]}^\dagger R_{[i_k j_k]} v_{[j_k]}}.$$

Formally, we regard  $\{i_1, \dots, i_p, j_1, \dots, j_p\}$  as the set of  $2p$  (index) variables that take values in  $\mathcal{I}_1$ . Let  $\mathcal{B}_p$  be the collection of all partitions of  $\{i_1, \dots, i_p, j_1, \dots, j_p\}$  such that  $i_k, j_k$  are not in the same block for all  $k = 1, \dots, p$ . For  $\Gamma \in \mathcal{B}_p$ , let  $n(\Gamma)$  be the number of its blocks and define a set of  $\mathcal{I}_1$ -valued variables as

$$L(\Gamma) := \{b_1, \dots, b_{n(\Gamma)}\}. \quad (5.8)$$

Now it is convenient to regard  $\Gamma$  as a symbol-to-symbol function

$$\Gamma : \{i_1, \dots, i_p, j_1, \dots, j_p\} \rightarrow L(\Gamma), \quad (5.9)$$

such that each  $\Gamma^{-1}(b_k)$  is a block of the partition. Then we can rewrite the sum as

$$\left| \sum_{i \neq j} u_{[i]}^\dagger R_{[ij]} v_{[j]} \right|^p = \sum_{\Gamma \in \mathcal{B}_p} \sum_{\substack{b_l \in \mathcal{I}_1, \\ l=1, \dots, n(\Gamma)}}^* \prod_{k=1}^{p/2} u_{[\Gamma(i_k)]}^\dagger R_{[\Gamma(i_k)\Gamma(j_k)]} v_{[\Gamma(j_k)]} \cdot \prod_{k=p/2+1}^p \overline{u_{[\Gamma(i_k)]}^\dagger R_{[\Gamma(i_k)\Gamma(j_k)]} v_{[\Gamma(j_k)]}}, \quad (5.10)$$

where  $\Sigma^*$  denote the summation subject to the condition that the values of  $b_1, \dots, b_n$  are ordered as  $b_1 < b_2 < \dots < b_n$ . We pick one term from the above summation and denote

$$\Delta(\Gamma) := \prod_{k=1}^{p/2} u_{[\Gamma(i_k)]}^\dagger R_{[\Gamma(i_k)\Gamma(j_k)]} v_{[\Gamma(j_k)]} \cdot \prod_{k=p/2+1}^p \overline{u_{[\Gamma(i_k)]}^\dagger R_{[\Gamma(i_k)\Gamma(j_k)]} v_{[\Gamma(j_k)]}}. \quad (5.11)$$

**Notations:** For any  $b_k \in L$ , we can define a corresponding  $\mathcal{I}_2$ -valued variable  $\bar{b}_k$  in the obvious way, and we denote

$$[L] := \{b_1, \dots, b_n, \bar{b}_1, \dots, \bar{b}_n\}. \quad (5.12)$$

For notational convenience, we will also use letters  $i, j, k, l$  to denote the symbols in  $L$ .

## 5.2 String and string operators

During the proof we will frequently use the following resolvent identities for rescaled matrix  $R$ . They follows immediately from Lemma 3.3.

**Lemma 5.3** (Resolvent identities for  $R_{[ij]}$  groups). *For  $k \notin J$  and  $i, j \in \mathcal{I}_1 \setminus J \cup \{k\}$ , we have*

$$R_{[ij]}^{[J]} = R_{[ij]}^{[Jk]} + R_{[ik]}^{[J]} \left( R_{[kk]}^{[J]} \right)^{-1} R_{[kj]}^{[J]}, \quad (5.13)$$

$$\left( R_{[ii]}^{[J]} \right)^{-1} = \left( R_{[ii]}^{[Jk]} \right)^{-1} - \left( R_{[ii]}^{[J]} \right)^{-1} R_{[ik]}^{[J]} \left( R_{[kk]}^{[J]} \right)^{-1} R_{[ki]}^{[J]} \left( R_{[ii]}^{[Jk]} \right)^{-1}, \quad (5.14)$$

$$\left( R_{[ii]}^{[J]} \right)^{-1} = w^{-1/2} H_{[ii]}^{[J]} - w^{-1} \sum_{l, l' \notin J \cup \{i\}} H_{[il]}^{[J]} R_{[ll']}^{[J]} H_{[li']}^{[J]}. \quad (5.15)$$

Furthermore, for  $i \neq j$  and  $L$  defined in (5.8), we have

$$R_{[ij]}^{[L \setminus \{ij\}]} = R_{[ii]}^{[L \setminus \{ij\}]} S_{[ij]} R_{[jj]}^{[L \setminus \{ij\}]}, \quad \text{with } S_{[ij]} = -w^{-1/2} H_{[ij]} + w^{-1} \sum_{k, l \notin L} H_{[ik]} R_{[kl]}^{[L]} H_{[lj]}. \quad (5.16)$$

In this section, we expand the  $R$  variables in  $\Delta(\Gamma)$  using the identities in Lemma 5.3. During the expansion, we need to distinguish carefully between an algebraic expression and its values as a random variable.

**Definition 5.4** (Strings). *Let  $\mathfrak{A}$  be an alphabet containing all symbols that may appear during the expansion, such as  $R_{[ij]}^{[J]}$ ,  $\left( R_{[ij]}^{[J]} \right)^{-1}$ ,  $S_{[ij]}$ ,  $u_{[i]}^\dagger$  and  $v_{[j]}$  for  $i, j, J \subset L(\Gamma)$ . We define a string  $\mathbf{s}$  to be a formal expression consisting of the symbols from  $\mathfrak{A}$ , and denote by  $\llbracket \mathbf{s} \rrbracket$  the random variable represented by it. Let  $\mathfrak{M}$  be the collection of all possible strings. We denote an empty string by  $\emptyset$ .*

Given a string  $\mathbf{s}$ , after an expansion of  $R$ 's in it, we will get a different string  $\mathbf{s}'$ . However they represent the same random variable  $\llbracket \mathbf{s} \rrbracket = \llbracket \mathbf{s}' \rrbracket$ . During the proof, we will identify more elements of  $\mathfrak{A}$  (see the symbols in (5.32)).

To perform the expansions in a systematical way, we define the following operators acting on strings. We call the symbols  $R_{[ij]}^{[J]}$ ,  $\left( R_{[ij]}^{[J]} \right)^{-1}$  to be *maximally expanded* if  $J \cup \{i, j\} = L$ . We call a string  $\mathbf{s}$  to be *maximally expanded* if all the  $R$  symbols in  $\mathbf{s}$  is maximally expanded.

**Definition 5.5** (String operators). (i) Define an operator  $\tau_0^{(k)}$  for  $\Omega \in \mathfrak{M}$ , in the following sense. Find the first  $R_{[ij]}^{[J]}$  in  $\Omega$  such that  $k \notin J \cup \{i, j\}$ , or the first  $\left(R_{[ii]}^{[J]}\right)^{-1}$  such that  $k \notin J \cup \{i\}$ . If  $R_{[ij]}^{[J]}$  is found, replace it with  $R_{[ij]}^{[Jk]}$ ; if  $\left(R_{[ii]}^{[J]}\right)^{-1}$  is found, replace it with  $\left(R_{[ii]}^{[Jk]}\right)^{-1}$ ; if neither is found,  $\tau_0^{(k)}(\Omega) = \Omega$  and we say that  $\tau_0^{(k)}$  is trivial for  $\Omega$ .

(ii) Define an operator  $\tau_1^{(k)}$  for  $\Omega \in \mathfrak{M}$ , in the following sense. Find the first  $R_{[ij]}^{[J]}$  in  $\Omega$  such that  $k \notin J \cup \{i, j\}$ , or the first  $\left(R_{[ii]}^{[J]}\right)^{-1}$  such that  $k \notin T \cup \{i\}$ . If  $R_{[ij]}^{[J]}$  is found, replace it with  $R_{[ik]}^{[J]} \left(R_{[kj]}^{[J]}\right)^{-1} R_{[kj]}^{[J]}$ ; if  $\left(R_{[ii]}^{[J]}\right)^{-1}$  is found, replace it with  $-\left(R_{[ii]}^{[J]}\right)^{-1} R_{[ik]}^{[J]} \left(R_{[kj]}^{[J]}\right)^{-1} R_{[ki]}^{[J]} \left(R_{[ii]}^{[Jk]}\right)^{-1}$ ; if neither is found,  $\tau_1^{(k)}(\Omega) = \emptyset$  and we say that  $\tau_1^{(k)}$  is null for  $\Omega$ .

(iii) Define an operator  $\rho$  for  $\Omega \in \mathfrak{M}$ , in the following sense. Find each maximally expanded  $R_{[ij]}^{[L \setminus \{ij\}]}$  in  $\Omega$  and replace it with  $R_{[ii]}^{[L \setminus \{ij\}]} S_{[ij]} R_{[jj]}^{[L \setminus \{ij\}]}$ . If nothing is found,  $\rho(\Omega) = \Omega$ .

According to Lemma 5.3, for any  $\Omega \in \mathfrak{M}$  we have

$$\left\| \left( \tau_0^{(k)} + \tau_1^{(k)} \right) (\Omega) \right\| = \|\Omega\|, \quad \|\rho(\Omega)\| = \|\Omega\| \quad (5.17)$$

**Definition 5.6.** Define the function  $\mathcal{F}_{d-\max} : \mathfrak{M} \rightarrow \mathbb{N}$  (where the subscript “d-max” stands for “distance to being maximally expanded”) through

$$\mathcal{F}_{d-\max} \left( R_{[ij]}^{[J]*} \right) = |L \setminus (J \cup \{i, j\})|,$$

where  $*$  could be 1 or  $-1$ , and

$$\mathcal{F}_{d-\max}(\Omega) = \sum_{R \text{ variables in } \Omega} \mathcal{F}_{d-\max}(R).$$

Define another function  $\mathcal{F}_{\text{off}} : \mathfrak{M} \rightarrow \mathbb{N}$  with  $\mathcal{F}_{\text{off}}(\Omega)$  being the number of off-diagonal symbols in  $\Omega$ .

By off-diagonal symbols, we mean the terms of the form  $A_{st}$  with  $s \notin \{t, \bar{t}\}$  or  $A_{[ij]}$  with  $i \neq j$ , e.g.  $R_{[ij]}^{[J]}$  and  $S_{[ij]}$  with  $i \neq j$ . Later we will define other types of off-diagonal symbols (see (5.32)). Note that a  $R$  symbol is maximally expanded if and only if  $\mathcal{F}_{d-\max}(R) = 0$  and a string  $\Omega$  is maximally expanded if and only if  $\mathcal{F}_{d-\max}(\Omega) = 0$ . The next two lemmas are almost trivial by Definition 5.5.

**Lemma 5.7.** If  $\tau_0^{(k)}(\Omega) = \Omega$  and  $\tau_1^{(k)}(\Omega) = \emptyset$ ,

$$\mathcal{F}_{d-\max} \left( \tau_0^{(k)}(\Omega) \right) = \mathcal{F}_{d-\max}(\Omega), \quad \mathcal{F}_{d-\max} \left( \tau_1^{(k)}(\Omega) \right) = 0; \quad (5.18)$$

otherwise,

$$\mathcal{F}_{d-\max} \left( \tau_0^{(k)}(\Omega) \right) = \mathcal{F}_{d-\max}(\Omega) - 1, \quad \mathcal{F}_{d-\max} \left( \tau_1^{(k)}(\Omega) \right) \leq \mathcal{F}_{d-\max}(\Omega) + 4n(\Gamma). \quad (5.19)$$

For  $\rho$ , we have

$$\mathcal{F}_{d-\max}(\rho(\Omega)) = \mathcal{F}_{d-\max}(\Omega) + a, \quad (5.20)$$

where  $a$  is the number of maximally expanded off-diagonal  $R$ 's in  $\Omega$ .

**Lemma 5.8.** For any  $\Omega \in \mathfrak{M}$ , we have

$$\mathcal{F}_{\text{off}} \left( \tau_0^{(k)}(\Omega) \right) = \mathcal{F}_{\text{off}}(\Omega), \quad \mathcal{F}_{\text{off}}(\rho(\Omega)) = \mathcal{F}_{\text{off}}(\Omega), \quad (5.21)$$

and

$$\mathcal{F}_{\text{off}}(\Omega) + 1 \leq \mathcal{F}_{\text{off}} \left( \tau_1^{(k)}(\Omega) \right) \leq \mathcal{F}_{\text{off}}(\Omega) + 2 \quad \text{if } \tau_1^{(k)}(\Omega) \neq \emptyset. \quad (5.22)$$

### 5.3 Expansion of the strings

For simplicity of notations, throughout the rest of this section we omit the complex conjugates on the right hand side of (5.11) (if we keep the complex conjugates, the proof is the same but with slightly heavier notations). Suppose the right hand side of (5.11) is represented by a string  $\Omega_\Delta$ . Given a binary word  $\mathbf{w} = a_1 a_2 \dots a_m$  with  $a_i \in \{0, 1\}$ , we define the operation

$$(\Omega_\Delta)_{\mathbf{w}} = \rho\tau_{a_m}^{(b_m)} \dots \rho\tau_{a_2}^{(b_2)} \rho\tau_{a_1}^{(b_1)} (\Omega_\Delta) \quad (5.23)$$

where  $b_{qn+r} := b_r$  (recall (5.8)) for any  $1 \leq r \leq n$  and  $q \in \mathbb{N}$ . So a binary words  $\mathbf{w}$  uniquely determines an operator composition. By (5.17),  $\llbracket (\Omega_\Delta)_{\mathbf{w}0} \rrbracket + \llbracket (\Omega_\Delta)_{\mathbf{w}1} \rrbracket = \llbracket (\Omega_\Delta)_{\mathbf{w}} \rrbracket$  and so we get

$$\sum_{|\mathbf{w}|=m} \llbracket (\Omega_\Delta)_{\mathbf{w}} \rrbracket = \llbracket \Omega_\Delta \rrbracket$$

for any  $m \geq 1$ , where  $|\mathbf{w}|$  is the length of  $\mathbf{w}$ .

**Lemma 5.9.** *Given any  $\mathbf{w}$  such that  $|\mathbf{w}| = (n^2 + 1)(p + 6l_0)$  and  $(\Omega_\Delta)_{\mathbf{w}} \neq \emptyset$ , either  $\mathcal{F}_{\text{off}}((\Omega_\Delta)_{\mathbf{w}}) \geq l_0 := (8/\zeta + 2)p$ , or  $(\Omega_\Delta)_{\mathbf{w}}$  is maximally expanded.*

*Proof.* We use  $m_0$  to denote the number of 0's in  $\mathbf{w}$ , and  $m_1$  to denote the number of 1's. Furthermore, we use  $m_0^{(0)}$  to denote the number of 0's corresponding to the trivial  $\tau_0$ 's, and  $m_0^{(1)}$  to denote the number of 0's corresponding to the non-trivial  $\tau_0$ 's. Assume  $\mathcal{F}_{\text{off}}((\Omega_\Delta)_{\mathbf{w}}) < l_0$  and  $(\Omega_\Delta)_{\mathbf{w}}$  is not maximally expanded. By (5.21)-(5.22),  $m_1 \leq l_0 - p \leq l_0$ . By (5.18)-(5.20),

$$\mathcal{F}_{\text{d-max}}((\Omega_\Delta)_{\mathbf{w}}) \leq \mathcal{F}_{\text{d-max}}(\Omega_\Delta) + l_0 + 4nm_1 - m_0^{(1)}.$$

Using  $\mathcal{F}_{\text{d-max}}(\Omega_\Delta) = np$ , we get a rough estimate  $m_0^{(1)} + m_1 < n(p + 6l_0)$ . By pigeonhole principle, there are at least  $n$  0's in a row in  $\mathbf{w}$  that correspond to trivial  $\tau_0$ 's. This indicates that  $(\Omega_\Delta)_{\mathbf{w}}$  is maximally expanded, which gives a contradiction.  $\square$

**Lemma 5.10.** *There exists constants  $C_{p,l_0}, C_{p,\zeta} > 0$  such that*

$$\sum_{\Gamma \in \mathcal{B}_p} \sum_{\substack{b_l \in \mathcal{I}_1, \\ l=1, \dots, n(\Gamma)}}^* \left| \mathbb{E} \sum_{\substack{|\mathbf{w}|=(n^2+1)(p+6l_0), \\ \mathcal{F}_{\text{off}}((\Omega_\Delta(\Gamma))_{\mathbf{w}}) \geq l_0}} \llbracket (\Omega_\Delta(\Gamma))_{\mathbf{w}} \rrbracket \right| \leq C_{p,l_0} N^{2p} \Phi^{l_0} \leq C_{p,\zeta} \Phi^p. \quad (5.24)$$

*Proof.* The first bound is due to the fact that each summand is bounded by  $C\Phi^{l_0}$  and there are at most  $N^{2p}$  of them. For the second bound, we used  $\Phi \leq CN^{-\zeta/2}$ .  $\square$

This lemma shows that all the strings with sufficiently many off-diagonal symbols contributes at most  $\Phi^p$ . It only remains to handle the maximally expanded strings. Define a diagonal symbol as

$$S_{[ii]} := - \begin{pmatrix} 0 & d_i X_{i\bar{i}} \\ \bar{d}_i X_{i\bar{i}}^\dagger & 0 \end{pmatrix} + w^{-1} \sum_{k, l \notin L} H_{[ik]} R_{[kl]}^{[L]} H_{[li]}, \quad (5.25)$$

such that

$$\left( R_{[ii]}^{[L \setminus \{i\}]} \right)^{-1} = \begin{pmatrix} -w^{1/2} & -z \\ -\bar{z} & -w^{1/2} \end{pmatrix} - S_{[ii]}. \quad (5.26)$$



Notice all the  $R$  symbols in a maximally expanded string is diagonal. We taylor expand  $R_{[ii]}^{[L \setminus \{i\}]}$  as

$$R_{[ii]}^{[L \setminus \{i\}]} = \left[ w^{-1/2} \pi_{[i]c}^{-1} + (S_{[ii]} - B_i) \right]^{-1} = \sum_{k=0}^{l_0-1} \tilde{\pi}_{ic} [(S_{[ii]} - B_i) \tilde{\pi}_{ic}]^k + O_{<}(\Phi^{l_0}), \quad (5.27)$$

where  $\tilde{\pi}_{[i]c} = w^{1/2} \pi_{[i]c}$ ,  $B_i = \begin{pmatrix} w^{1/2} |d_i|^2 m_{2c} & 0 \\ 0 & w^{1/2} m_{1c} \end{pmatrix}$ , and for the error term,

$$S_{[ii]} - B_i = w^{-1/2} Z_{[i]}^{[L \setminus \{i\}]} + w^{1/2} \begin{pmatrix} |d_i|^2 (m_{2c} - m_2^{[L]}) & 0 \\ 0 & m_{1c} - m_1^{[L]} \end{pmatrix} < \Phi$$

by (4.15) and the averaged local law. Now for all maximally expanded  $(\Omega_\Delta)_{\mathbf{w}}$  with  $|\mathbf{w}| = (n^2 + 1)(p + 6l_0)$ , denote by  $\sigma[(\Omega_\Delta)_{\mathbf{w}}]$  the expression after plugging in (5.26) and (5.27) without the tail terms. Similar to Lemma 5.10, we have

$$\sum_{\Gamma \in \mathcal{B}_p} \sum_{\substack{b_l \in \mathcal{I}_1, \\ l=1, \dots, n(\Gamma)}}^* \left| \mathbb{E} \sum_{\substack{|\mathbf{w}|=(n^2+1)(p+6l_0), \\ (\Omega_\Delta)_{\mathbf{w}} \text{ maximally expanded}}} \left( [(\Omega_{\Delta(\Gamma)})_{\mathbf{w}}] - \sigma[(\Omega_{\Delta(\Gamma)})_{\mathbf{w}}] \right) \right| \leq C_{p,\zeta} \Phi^p.$$

From the above bound and Lemmas 5.9, 5.10, we see that to prove (5.7), it suffices to show

$$\sum_{\Gamma \in \mathcal{B}_p} \sum_{\substack{b_l \in \mathcal{I}_1, \\ l=1, \dots, n(\Gamma)}}^* \left| \mathbb{E} \sum_{\substack{|\mathbf{w}|=(n^2+1)(p+6l_0), \\ (\Omega_\Delta)_{\mathbf{w}} \text{ maximally expanded}}} \sigma[(\Omega_{\Delta(\Gamma)})_{\mathbf{w}}] \right| \leq C_{p,\zeta} \Phi^p. \quad (5.28)$$

We write  $\sigma[(\Omega_\Delta)_{\mathbf{w}}]$  as a sum of monomials in terms of  $S_{[ij]}$ ,

$$\sigma[(\Omega_\Delta)_{\mathbf{w}}] = \sum_i M(\mathbf{w}, \Delta(\Gamma), i), \quad (5.29)$$

where  $i$  is an index to label these monomials. Notice that after plugging (5.29) into (5.28), the number of summands  $M(\mathbf{w}, \Delta(\Gamma), i)$  inside the expectation only depends on  $p$  and  $\zeta$ . Thus to show (5.28), it suffices to prove the following lemma.

**Lemma 5.11.** *Fix any  $\Gamma \in \mathcal{B}_p$  and binary word  $\mathbf{w}$  with  $|\mathbf{w}| = (n^2 + 1)(p + 6l_0)$ . Suppose  $(\Omega_\Delta)_{\mathbf{w}}$  is maximally expanded. Let  $M(\mathbf{w}, \Delta(\Gamma))$  be a monomial in  $\sigma[(\Omega_{\Delta(\Gamma)})_{\mathbf{w}}]$ . We have*

$$\sum_{b_l \in \mathcal{I}_1, l=1, \dots, n(\Gamma)}^* |\mathbb{E} M(\mathbf{w}, \Delta(\Gamma))| \leq C_{p,\zeta} \Phi^p \quad (5.30)$$

for some constant  $C_{p,\zeta}$  that only depends on  $p$  and  $\zeta$ .

For the rest of this section, we fix a  $\Gamma \in \mathcal{B}_p$  and a maximally expanded  $(\Omega_{\Delta(\Gamma)})_{\mathbf{w}}$  with  $|\mathbf{w}| = (n^2 + 1)(p + 6l_0)$ . Then we fix a monomial  $M(\mathbf{w}, \Delta(\Gamma))$  in  $\sigma[(\Omega_{\Delta(\Gamma)})_{\mathbf{w}}]$ . Let  $\Omega_M$  be the string form of  $M(\mathbf{w}, \Delta(\Gamma))$  in terms of  $S_{[ij]}$ . It is not hard to see that

$$\mathcal{F}_{\text{off}}(\Omega_M) = \mathcal{F}_{\text{off}}((\Omega_\Delta)_{\mathbf{w}}). \quad (5.31)$$

Now we decompose  $S_{[ij]}$  as

$$S_{[ij]} = S_{ij}^X + S_{ij}^X + S_{ij}^R + S_{ij}^R + S_{ij}^R + S_{ij}^R, \quad (5.32)$$

where we define the following symbols in  $\mathfrak{A}$ :

$$S_{ij}^X := d_i X_{i\bar{j}} \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}, \quad S_{ij}^X := \bar{d}_i X_{ij}^\dagger \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}, \quad (5.33)$$

$$S_{ij}^R := \sum_{k,l \notin L} d_i d_l X_{i\bar{k}} X_{l\bar{j}} \begin{pmatrix} 0 & R_{kl}^{[L]} \\ 0 & 0 \end{pmatrix}, \quad S_{ij}^R := \sum_{k,l \notin L} d_i \bar{d}_l X_{i\bar{k}} X_{l\bar{j}}^\dagger \begin{pmatrix} R_{kl}^{[L]} & 0 \\ 0 & 0 \end{pmatrix}, \quad (5.34)$$

$$S_{ij}^R := \sum_{k,l \notin L} \bar{d}_i d_l X_{ik}^\dagger X_{l\bar{j}} \begin{pmatrix} 0 & 0 \\ R_{kl}^{[L]} & 0 \end{pmatrix}, \quad S_{ij}^R := \sum_{k,l \notin L} \bar{d}_i \bar{d}_l X_{ik}^\dagger X_{l\bar{j}}^\dagger \begin{pmatrix} 0 & 0 \\ R_{kl}^{[L]} & 0 \end{pmatrix}. \quad (5.35)$$

We expand  $S_{[ij]}$ 's of  $M(\mathbf{w}, \Delta(\Gamma))$  as in (5.32), and write  $M(\mathbf{w}, \Delta(\Gamma))$  as a sum of monomials in terms of  $S_{st}^X$  and  $S_{st}^R$ ,

$$M(\mathbf{w}, \Delta(\Gamma)) = \sum_i Q(\mathbf{w}, \Delta(\Gamma), i), \quad (5.36)$$

where  $i$  is an index to label these monomials. Again it is not hard to see that

$$\mathcal{F}_{\text{off}}(\Omega_Q) = \mathcal{F}_{\text{off}}(\Omega_M) = \mathcal{F}_{\text{off}}((\Omega_\Delta)_\mathbf{w}). \quad (5.37)$$

Since the number of summands in (5.36) is independent of  $N$ , to prove (5.30) it suffices to show

$$\sum_{b_l \in \mathcal{I}_1, l=1, \dots, n(\Gamma)}^* |\mathbb{E} Q(\mathbf{w}, \Delta(\Gamma))| \leq C_{p,\zeta} \Phi^p \quad (5.38)$$

for any monomial  $Q(\mathbf{w}, \Delta(\Gamma))$  in (5.36). Throughout the following, we fix a  $Q(\mathbf{w}, \Delta(\Gamma))$  with nonzero expectation, and denote by  $\Omega_Q$  the string form of  $Q(\mathbf{w}, \Delta(\Gamma))$  in terms of  $S_{st}^X$  and  $S_{st}^R$ . Notice the  $R$  variables in  $S_{st}^R$  are maximally expanded. As a result, the  $S_{st}^X$  variables are independent of  $S_{st}^R$  variables in  $Q(\mathbf{w}, \Delta(\Gamma))$ . Therefore we make the following observation: if  $S_{st}^X$  appears as a symbol in  $\Omega_Q$ , then  $\Omega_Q$  contains at least two of them.

**Definition 5.12.** Recall  $\Gamma$  defined in (5.9). Let  $h$  be the number of blocks of  $\Gamma$  whose size is 1, i.e.

$$h := \sum_{l=1}^{n(\Gamma)} \mathbf{1}(|\Gamma^{-1}(b_l)| = 1). \quad (5.39)$$

For  $l = 1, \dots, n$ , define

$$I_l := |\{i_1, \dots, i_p\} \cap \Gamma^{-1}(b_l)|, \quad J_l := |\{j_1, \dots, j_p\} \cap \Gamma^{-1}(b_l)|.$$

**Lemma 5.13.** Suppose for any  $b_1, \dots, b_n$  taking distinct values in  $\mathcal{I}_1$ ,

$$|\mathbb{E} Q(\mathbf{w}, \Delta(\Gamma))| \leq C N^{-h/2} \Phi^p \prod_{l=1}^n |u_{[b_l]}|^{I_l} |v_{[b_l]}|^{J_l} \quad (5.40)$$

holds for some constant  $C$  independent of  $N$ . Then the estimate (5.38) holds.

*Proof.* By Cauchy-Schwarz inequality,

$$\sum_{k=1}^N |u_{[k]}|^a |v_{[k]}|^b \leq \begin{cases} N^{1/2} & \text{if } a + b = 1 \\ 1 & \text{if } a + b \geq 2 \end{cases}.$$

Then using  $h = \sum_{l=1}^n \mathbf{1}(I_l + J_l = 1)$ , we get

$$\sum_{b_l \in \mathcal{I}_1, l=1, \dots, n(\Gamma)}^* |\mathbb{E}Q(\mathbf{w}, \Delta(\Gamma))| \leq C\Phi^p N^{-h/2} \prod_{l=1}^n \sum_{b_l \in \mathcal{I}_1} |u_{[b_l]}|^{I_l} |v_{[b_l]}|^{J_l} \leq C\Phi^p.$$

□

Hence it suffices to prove (5.40). The key is to extract the  $N^{-h/2}$  factor from  $\mathbb{E}Q(\mathbf{w}, \Delta(\Gamma))$ . For this purpose, we need to keep track of the indices in  $L$  during the expansion.

**Definition 5.14.** Define a function  $\mathcal{F}_{\text{in}} : L \times \mathfrak{M} \rightarrow \mathbb{N}$  with  $\mathcal{F}_{\text{in}}(l, \Omega)$  giving the number of times  $l$  or  $\bar{l}$  appears as an index of off-diagonal  $R$  or  $S$  in  $\Omega$ .

The following lemma follows immediately from Definition 5.5 and the expansions we have done to obtain  $\Omega_Q$  from  $(\Omega_\Delta)_{\mathbf{w}}$ .

**Lemma 5.15.** (1) For any string  $\Omega$ , if  $\tau_0^{(k)}$  is not trivial for  $\Omega$ , then

$$\mathcal{F}_{\text{in}}(l, \tau_0^{(k)}(\Omega)) = \mathcal{F}_{\text{in}}(l, \Omega), \quad \mathcal{F}_{\text{in}}(l, \tau_1^{(k)}(\Omega)) = \mathcal{F}_{\text{in}}(l, \Omega) + 2\delta_{kl}. \quad (5.41)$$

(2) For any string  $\Omega$ ,

$$\mathcal{F}_{\text{in}}(l, \rho(\Omega)) = \mathcal{F}_{\text{in}}(l, \Omega). \quad (5.42)$$

(3) For any maximally expanded  $(\Omega_\Delta)_{\mathbf{w}}$ ,

$$\mathcal{F}_{\text{in}}(l, \Omega_Q) = \mathcal{F}_{\text{in}}(l, (\Omega_\Delta)_{\mathbf{w}}). \quad (5.43)$$

Let  $\Omega_Q^X$  be the substring of  $\Omega_Q$  containing only  $S^X$  symbols, and  $\Omega_Q^R$  be the substring of  $\Omega_Q$  containing only  $S^R$  symbols. Define

$$\mathcal{V} := \{l \in L \mid \mathcal{F}_{\text{in}}(l, \Omega_\Delta) = 1\}, \quad (5.44)$$

and

$$\mathcal{V}_0 := \{l \in L \mid \mathcal{F}_{\text{in}}(l, \Omega_\Delta) = 1 \text{ and } \mathcal{F}_{\text{in}}(l, \Omega_Q^X) = 0\}, \quad (5.45)$$

$$\mathcal{V}_1 := \{l \in L \mid \mathcal{F}_{\text{in}}(l, \Omega_\Delta) = 1 \text{ and } \mathcal{F}_{\text{in}}(l, \Omega_Q^X) \geq 2\}. \quad (5.46)$$

Recall the observation above Definition 5.12,  $\mathcal{V} = \mathcal{V}_0 \cup \mathcal{V}_1$  and

$$h = |\mathcal{V}| = |\mathcal{V}_0| + |\mathcal{V}_1|.$$

Let  $n_X$  be the number of off-diagonal  $S^X$  symbols in  $\Omega_Q^X$  and  $n_R$  be the number of off-diagonal  $S^R$  symbols in  $\Omega_Q^R$ . Notice that  $n_o := n_X + n_R$  is the total number of off-diagonal symbols in  $\Omega_Q$ .

## 5.4 Introduction of graphs and conclusion of the proof

We introduce the graphs to conclude the proof of (5.40). We use a connected graph to represent the string  $\Omega_Q$ , call it by  $\mathfrak{G}_{Q0}$ . The indices in  $[L]$  are represented by black nodes in  $\mathfrak{G}_{Q0}$ . The  $S_{st}^X$  or  $S_{st}^R$  symbols in  $\Omega_Q$  are represented by edges connecting the nodes  $s$  and  $t$ . We also define colors for the nodes and edges, where the color set for nodes is  $\{black, white\}$  and the color set for edges is  $\{S^X, S^R, X, R\}$ . In  $\mathfrak{G}_{Q0}$ , all the nodes are black, all  $S^X$  edges are assigned  $S^X$  color and all  $S^R$  edges are assigned  $S^R$  color. We show a possible graph in Fig. 3. In this subsection, we identify an index with its node representation, and a symbol with its edge representation.

**Definition 5.16.** Define function  $\deg$  on the nodes set  $[L]$ , where  $\deg(l)$  is the number of  $S^R$  edges connecting to the node  $l$ .

By Lemma 5.15, we see that for any  $l \in \mathcal{V}_0$ ,

$$\mathcal{F}_{\text{in}}(l, \Omega_Q) \equiv \deg(l) + \deg(\bar{l}) \equiv 1 \pmod{2}. \quad (5.47)$$

Hence

$$|\mathcal{V}_0| = \sum_{l \in \mathcal{V}_0} [\mathcal{F}_{\text{in}}(l, \Omega_Q) \pmod{2}] \leq \sum_{l \in \mathcal{V}_0} [(\deg(l) \pmod{2}) + (\deg(\bar{l}) \pmod{2})]. \quad (5.48)$$

Now we expand the  $S^R$  edges. Take the  $S_{ij}^R$  edge as an example (recall (5.34)). We replace the  $S_{ij}^R$  edge with an  $R$ -group, defined as following. We add two white colored nodes to represent the summation indices  $\bar{k}, l \notin [L]$ , two  $X$ -colored edges to represent  $X_{i\bar{k}}$  and  $X_{l\bar{j}}$ , and a  $R$ -colored edge connecting  $\bar{k}$  and  $l$  to represent  $\begin{pmatrix} 0 & R_{\bar{k}l}^{[L]} \\ 0 & 0 \end{pmatrix}$ . We call the subgraph consisting of the three new edges and their nodes an  $R$ -group. If  $i = j$ , we call it a diagonal  $R$ -group; otherwise, call it an off-diagonal  $R$ -group. We expand all  $S^R$  edges in  $\mathfrak{G}_{Q0}$  into  $R$ -groups and call the resulting graph  $\mathfrak{G}_{Q1}$ . For example, after expanding the  $S^R$  edges in Fig.3, we get the graph in Fig.4. In the graph  $\mathfrak{G}_{Q1}$ , the  $R$  edges,  $X$  edges and  $S^X$  edges are mutually independent, since the  $R$  symbols are maximally expanded, and the white nodes are different from the black nodes.

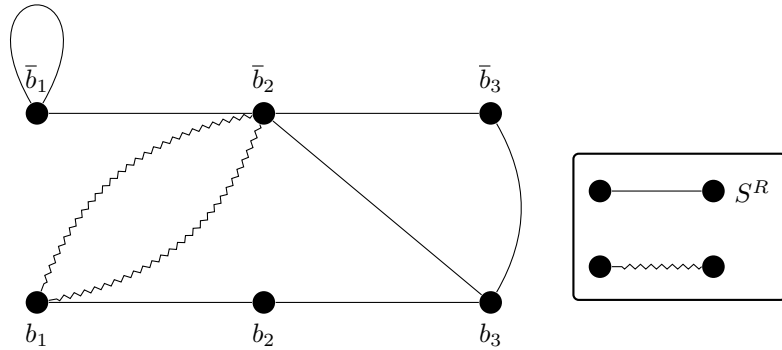


Figure 3: An example of the graph  $\mathfrak{G}_{Q0}$ .

Notice that each white node represents a summation index. As we have done for the black nodes, we first partition the white nodes into blocks and then assign values to the blocks when doing the

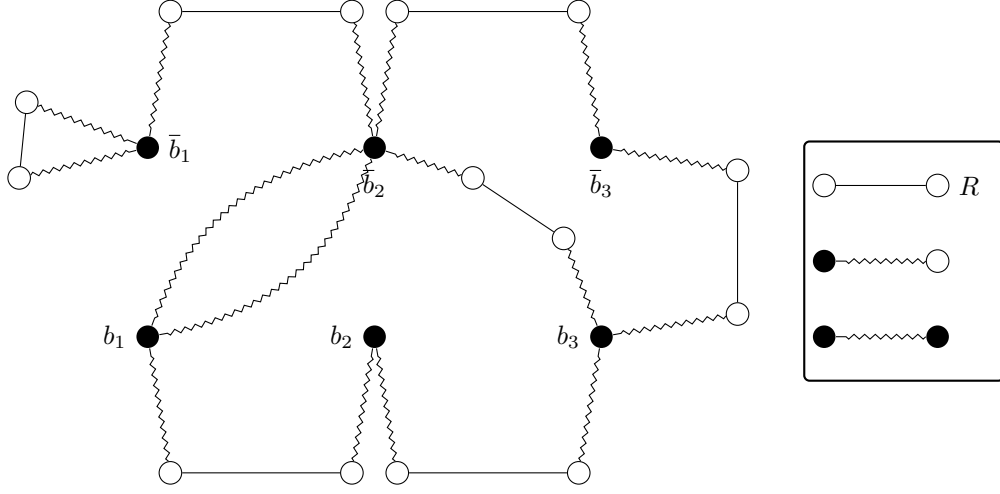


Figure 4: The resulting graph  $\mathfrak{G}_{Q1}$  after expanding each  $S^R$  in Fig. 3 into  $R$ -groups.

summation. Let  $W$  be the set of all white nodes in  $\mathfrak{G}_{Q1}$ , and let  $\mathcal{W}$  be the collection of all partitions of  $W$ . Fix a partition  $\gamma \in \mathcal{W}$  and denote its blocks by  $W_1, \dots, W_{m(\gamma)}$ . If two white nodes of some off-diagonal  $R$ -group happen to lie in the same block, then we merge the two nodes into one diamond white node (Fig. 5a). All the other white nodes are called normal (Fig. 5b). Let  $n_R^{(d)}$  be the number of diamond nodes ( $\leq$  the number of diagonal  $R$ -edges in  $\mathfrak{G}_{Q1}$ ). Then we trivially have

$$\# \text{ of white nodes} = -n_R^{(d)} + \sum_{k=1}^n [\deg(b_k) + \deg(\bar{b}_k)]. \quad (5.49)$$

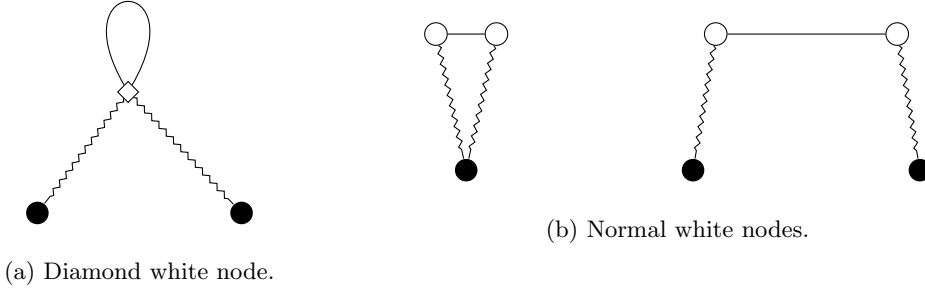


Figure 5: Two types of white nodes

By (5.48), there are  $|\mathcal{V}_0|$  black nodes with odd deg in  $[\mathcal{V}_0]$  (where  $[\mathcal{V}_0]$  is defined in the obvious way). WLOG, we assume these nodes are  $b_1, \dots, b_{|\mathcal{V}_0|}$ . To have nonzero expectation, each white block must contain at least two white nodes. Therefore for each  $k = 1, \dots, |\mathcal{V}_0|$ , there exists a block connecting to  $b_k$  which contains at least 3 white nodes. Call such a block  $W(b_k)$ , and denote by  $A(b_k)$  the set of the adjacent white nodes to  $b_k$  in  $W(b_k)$ . (Note that the  $W(b_k)$ 's or  $A(b_k)$ 's are not necessarily distinct.) WLOG, let  $W_1, \dots, W_d$  be the distinct blocks among all  $W(b_k)$ 's. Define

$$\mathcal{V}_{00} := \{b_k \mid A(b_k) \text{ has no normal white nodes}, 1 \leq k \leq |\mathcal{V}_0|\},$$

and

$$\mathcal{V}_{01} := \{b_k \mid A(b_k) \text{ has at least one normal white node}, 1 \leq k \leq |\mathcal{V}_0|\}.$$

The following lemma gives the key estimates we need.

**Lemma 5.17.** *For any partition  $\gamma \in \mathcal{W}$ ,*

$$m(\gamma) \leq \frac{-|\mathcal{V}_{00}| - |\mathcal{V}_{01}|/2 - n_R^{(d)} + \sum_{k=1}^n [\deg(b_k) + \deg(\bar{b}_k)]}{2}, \quad (5.50)$$

and

$$n_X + n_R \geq p + |\mathcal{V}_1| + |\mathcal{V}_{00}|, \quad n_X \geq |\mathcal{V}_1|, \quad n_R^{(d)} \geq |\mathcal{V}_{00}|. \quad (5.51)$$

*Proof.* The second inequality of (5.51) can be proved easily through

$$|\mathcal{V}_1| \leq |\{k \in L \mid \mathcal{F}_{\text{in}}(k, \Omega_Q^X) \geq 2\}| \leq n_X.$$

Notice for  $b_k \in \mathcal{V}_0$ ,  $A(b_k)$  contains at least three diamond white nodes, while each of the white node is share by another  $b_l$ . Thus we trivially have  $|\mathcal{V}_{00}| \leq n_R^{(d)}$ .

Now we prove (5.50). A diamond white node is connected to two black nodes and a normal white node is connected to one black node. Hence a diamond white node belongs to two sets  $A(b_{k_1}), A(b_{k_2})$ , and a normal white node belongs to exactly one set  $A(b_k)$ . Therefore for each  $i = 1, \dots, d$ , if  $W_i$  contains exactly one  $A(b_k)$  then

$$|W_i| \geq 3 \geq 2 + \mathbf{1}_{\mathcal{V}_{01}}(b_k) + \frac{\mathbf{1}_{\mathcal{V}_{00}}(b_k)}{2}.$$

Otherwise if  $W_i$  contains more than one  $A(b_k)$ , then

$$|W_i| \geq \sum_{b_k: A(b_k) \subseteq W_i} \left( 2 \cdot \mathbf{1}_{\mathcal{V}_{01}}(b_k) + \frac{3}{2} \cdot \mathbf{1}_{\mathcal{V}_{00}}(b_k) \right) \geq 2 + \sum_{b_k: A(b_k) \subseteq W_i} \left( \mathbf{1}_{\mathcal{V}_{01}}(b_k) + \frac{\mathbf{1}_{\mathcal{V}_{00}}(b_k)}{2} \right).$$

Here the first inequality can be understood as following. For each black node  $b_k$  with  $A(b_k) \subseteq W_i$ , we count the number of white nodes in  $A(b_k)$  and add them together. During the counting, we assign weight-1 to a normal white node and weight-1/2 to a diamond white node (since it is shared by two different black nodes). If  $b_k \in \mathcal{V}_{00}$ , there are at least three diamond white nodes in  $A(b_k)$  with total weight  $\geq 3/2$ ; if  $b_k \in \mathcal{V}_{01}$ , there are at least one normal white node and two other white nodes in  $A(b_k)$  with total weight  $\geq 2$ . Thus  $\sum_{b_k: A(b_k) \subseteq W_i} (2 \cdot \mathbf{1}_{\mathcal{V}_{01}}(b_k) + \frac{3}{2} \cdot \mathbf{1}_{\mathcal{V}_{00}}(b_k))$  is smaller than the number of white nodes in  $W_i$ . Then summing  $|W_i|$  over  $i$ , we get

$$\sum_{i=1}^d |W_i| \geq 2d + |\mathcal{V}_{01}| + \frac{|\mathcal{V}_{00}|}{2}.$$

For the other  $m - d$  blocks, each of them contains at least two white nodes. Therefore

$$2m + |\mathcal{V}_{01}| + \frac{|\mathcal{V}_{00}|}{2} \leq \sum_{i=1}^d |W_i| + 2(m - d) \leq -n_R^{(d)} + \sum_{k=1}^n [\deg(b_k) + \deg(\bar{b}_k)],$$

where we use (5.49) in the last step. This proves (5.50).

For  $b_k \in \mathcal{V}_{00}$ ,  $A(b_k)$  contains at least three white nodes from off-diagonal  $R$ -groups,

$$\mathcal{V}_{00} \subseteq \{b_k \in L \mid \mathcal{F}_{\text{in}}(b_k, \Omega_\Delta) = 1 \text{ and } \mathcal{F}_{\text{in}}(b_k, \Omega_Q^R) \geq 3\} =: \mathcal{V}_2.$$

Recall (5.41)-(5.42), only  $\tau_1^{(k)}$  may increase  $\mathcal{F}_{\text{in}}$ . Thus  $\mathbf{w}$  contains  $\tau_1^{(b_k)}$  for each  $b_k \in \mathcal{V}_1 \cup \mathcal{V}_2$  (recall the definition of  $\mathcal{V}_1$  in (5.46)). Therefore by (5.22), (5.37) and the fact that  $\mathcal{V}_{00}$  and  $\mathcal{V}_1$  are disjoint,

$$n_X + n_R = \mathcal{F}_{\text{off}}((\Omega_\Delta)_{\mathbf{w}}) \geq \mathcal{F}_{\text{off}}(\Omega_\Delta) + |\mathcal{V}_1 \cup \mathcal{V}_2| \geq p + |\mathcal{V}_1| + |\mathcal{V}_{00}|.$$

This proves the first inequality of (5.51).  $\square$

By (2.3) and (5.6), a diagonal  $R$  edge contributes 1, an off-diagonal  $R$  edge contributes  $\Phi$ , and  $S^X$  or  $X$  edge contributes  $N^{-1/2}$ . Denote

$$\mathcal{U} = \prod_{l=1}^n |u_{[b_l]}|^{I_l} |v_{[b_l]}|^{J_l}.$$

Then using Lemma 2.21, we get

$$\begin{aligned} |\mathbb{E}Q(\mathbf{w}, \Delta(\Gamma))| &\leq C\mathcal{U} \left(N^{-1/2}\right)^{n_X} \sum_{\gamma \in \mathcal{W}} \sum_{\gamma(W_1), \dots, \gamma(W_m) \in \mathcal{I} \setminus L}^* \Phi^{n_R - n_R^{(d)}} \prod_{k=1}^n \left(N^{-1/2}\right)^{\deg(b_k) + \deg(\bar{b}_k)} \\ &\leq C\mathcal{U} N^{-n_X/2} \sum_{\gamma \in \mathcal{W}} N^{m - \frac{\sum_{k=1}^n \deg(b_k) + \deg(\bar{b}_k)}{2}} \Phi^{n_R - n_R^{(d)}} \\ &\leq C\mathcal{U} N^{-n_X/2} \sum_{\gamma \in \mathcal{W}} N^{\frac{-|\mathcal{V}_{01}| - |\mathcal{V}_{00}|/2 - n_R^{(d)}}{2}} \Phi^{n_R - n_R^{(d)}} \\ &\leq C\mathcal{U} N^{-h/2} \sum_{\gamma \in \mathcal{W}} N^{-(n_X - |\mathcal{V}_1|)/2} N^{-(n_R^{(d)} - |\mathcal{V}_{00}|)/2} \Phi^{n_R - n_R^{(d)}} \\ &\leq C\mathcal{U} N^{-h/2} \sum_{\gamma \in \mathcal{W}} \Phi^{n_X + n_R - |\mathcal{V}_1| - |\mathcal{V}_{00}|} \leq C\mathcal{U} N^{-h/2} \Phi^p, \end{aligned}$$

where in the third step we used (5.50), in the fourth step  $h = |\mathcal{V}| = |\mathcal{V}_1| + |\mathcal{V}_{00}| + |\mathcal{V}_{01}|$ , in the fifth step  $N^{-1/2} \leq \Phi$  and (5.51), and in the last step (5.51). Thus we have proved (5.40), which concludes the proof of Proposition 5.1.

## 6 Anisotropic local law: self-consistent comparison

In this section we prove Theorem 2.19. We first prove the anisotropic and averaged local laws under the vanishing third moment assumption (2.23). When  $\eta \geq N^{-1/2+\zeta} |m_{2c}|^{-1}$ , the anisotropic and averaged local laws can be established without assuming (2.23). For convenience, we only consider the case  $w \in \mathbf{D}$  and  $|z|^2 \leq 1 - \tau$  in this section. The proof for other cases is almost the same.

Following the notations in the arguments between Theorems 2.18 and 2.19,

$$H(TX - z, w) = \bar{T} \begin{pmatrix} -w(D^\dagger D)^{-1} & w^{1/2}(V_1 X - (UD)^{-1}z) \\ w^{1/2}(V_1 X - (UD)^{-1}z)^\dagger & -wI \end{pmatrix} \bar{T}^\dagger, \quad \bar{T} := \begin{pmatrix} UD & 0 \\ 0 & I \end{pmatrix}. \quad (6.1)$$

Now we define

$$\mathcal{G}(w) := |w|^{1/2} \begin{pmatrix} -w(D^\dagger D)^{-1} & w^{1/2}(V_1 X - (UD)^{-1}z) \\ w^{1/2}(V_1 X - (UD)^{-1}z)^\dagger & -wI \end{pmatrix}^{-1} = |w|^{1/2} \bar{T}^\dagger G \bar{T}. \quad (6.2)$$

Since  $T$  is invertible and  $\|T\| + \|T^{-1}\| \leq \tau^{-1}$  by (2.4), to prove the anisotropic law in Theorem 2.19, it suffices to show

$$\|\mathcal{G}(w) - \tilde{\Pi}(w)\| < \Phi(w) \quad (6.3)$$

where

$$\tilde{\Pi}(w) := |w|^{1/2} \bar{T}^\dagger \Pi(w) \bar{T}, \quad \Phi(w) := |w|^{1/2} \Psi(w). \quad (6.4)$$

Notice we have  $\|\tilde{\Pi}\| = O(1)$  by (3.31). By the remark around (2.50), if  $X = X^{Gauss}$  is Gaussian, then (6.3) holds. Hence for a general  $X$ , it suffices to prove that

$$\|\mathcal{G}(X, w) - \mathcal{G}(X^{Gauss}, w)\| < \Phi(w). \quad (6.5)$$

Similar to Lemma 3.5, it is easy to prove the following estimates for  $\mathcal{G}$ .

**Lemma 6.1.** *For  $i \in \mathcal{I}_1^M$ , we define  $\mathbf{v}_i = V_1 \mathbf{e}_i \in \mathbb{R}^{\mathcal{I}_1}$ , i.e.  $\mathbf{v}_i$  is the  $i$ -th column vector of  $V_1$ . Let  $\mathbf{u} \in \mathbb{R}^{\mathcal{I}_1}$  and  $\mathbf{w} \in \mathbb{R}^{\mathcal{I}_2}$ , then we have for some constant  $C > 0$ ,*

$$\sum_{\mu \in \mathcal{I}_2} |\mathcal{G}_{\mathbf{w}\mu}|^2 = |w|^{1/2} \frac{\text{Im } \mathcal{G}_{\mathbf{w}\mathbf{w}}}{\eta}, \quad (6.6)$$

$$\sum_{i \in \mathcal{I}_1^M} |\mathcal{G}_{\mathbf{u}\mathbf{v}_i}|^2 \leq C |w|^{1/2} \frac{\text{Im } \mathcal{G}_{\mathbf{u}\mathbf{u}}}{\eta}, \quad (6.7)$$

$$\sum_{i \in \mathcal{I}_1^M} |\mathcal{G}_{\mathbf{w}\mathbf{v}_i}|^2 \leq C \left( |w|^{-1/2} \mathcal{G}_{\mathbf{w}\mathbf{w}} + \bar{w} |w|^{-1/2} \frac{\text{Im } \mathcal{G}_{\mathbf{w}\mathbf{w}}}{\eta} \right), \quad (6.8)$$

$$\sum_{\mu \in \mathcal{I}_2} |\mathcal{G}_{\mathbf{u}\mu}|^2 \leq C \left( |w|^{-1/2} \mathcal{G}_{\mathbf{u}\mathbf{u}} + \bar{w} |w|^{-1/2} \frac{\text{Im } \mathcal{G}_{\mathbf{u}\mathbf{u}}}{\eta} \right), \quad (6.9)$$

## 6.1 Self-consistent comparison

Our proof basically follows the arguments in [24, Section 7] with some minor modifications. Thus we will not write down all the details for the proof. By polarization, it suffices to show the following proposition.

**Proposition 6.2.** *Fix  $|z|^2 \leq 1 - \tau$  and suppose that the assumptions of Theorem 2.19 hold. If (2.23) holds or  $\eta \geq N^{-1/2+\zeta} |m_{2c}|^{-1}$ , then for any regular domain  $\mathbf{S} \subseteq \mathbf{D}$ ,*

$$\left\langle \mathbf{v}, \left( \mathcal{G}(w) - \tilde{\Pi}(w) \right) \mathbf{v} \right\rangle < \Phi(w) \quad (6.10)$$

uniformly in  $w \in \mathbf{S}$  and any deterministic unit vectors  $\mathbf{v} \in \mathbb{C}^{\mathcal{I}}$ .

We first assume that (2.23) holds. Then we will show how to modify the arguments to prove the  $\eta \geq N^{-1/2+\zeta} |m_{2c}|^{-1}$  case. The proof consists of a bootstrap argument from larger scales to smaller scales in multiplicative increments of  $N^{-\delta}$ , where

$$\delta \in \left( 0, \frac{\zeta}{2C_0} \right), \quad (6.11)$$

with  $C_0 > 0$  being a universal constant that will be chosen large enough in the proof. For any  $\eta \geq |m_{1c}|^{-1} N^{-1+\zeta}$ , we define

$$\eta_l := \eta N^{\delta l} \text{ for } l = 0, \dots, L-1, \quad \eta_L := 1. \quad (6.12)$$



where  $L \equiv L(\eta) := \max \{l \in \mathbb{N} \mid \eta N^{\delta(l-1)} < 1\}$ . Note that  $L \leq 2\delta^{-1}$ .

By (3.13), the function  $w \mapsto \mathcal{G}(w) - \tilde{\Pi}(w)$  is Lipschitz continuous in  $\mathbf{S}$  with Lipschitz constant bounded by  $CN^3$ . Thus to prove (6.10) for all  $w \in \mathbf{S}$ , it suffices to show (6.10) holds for all  $w$  in some discrete but sufficiently dense subset  $\hat{\mathbf{S}} \subset \mathbf{S}$ . We will use the following discretized domain  $\hat{\mathbf{S}}$ .

**Definition 6.3.** Let  $\hat{\mathbf{S}}$  be an  $N^{-10}$ -net of  $\mathbf{S}$  such that  $|\hat{\mathbf{S}}| \leq N^{20}$  and

$$E + i\eta \in \hat{\mathbf{S}} \Rightarrow E + i\eta_l \in \hat{\mathbf{S}} \text{ for } l = 1, \dots, L(\eta).$$

The bootstrapping is formulated in terms of two scale-dependent properties  $(\mathbf{A}_m)$  and  $(\mathbf{C}_m)$  defined on the subsets

$$\hat{\mathbf{S}}_m := \left\{ w \in \hat{\mathbf{S}} \mid \operatorname{Im} w \geq N^{-\delta m} \right\}.$$

$(\mathbf{A}_m)$  For all  $w \in \hat{\mathbf{S}}_m$ , all deterministic unit vector  $\mathbf{v}$ , and all  $X$  satisfying (2.2)-(2.3), we have

$$\operatorname{Im} \mathcal{G}_{\mathbf{v}\mathbf{v}}(w) < |w|^{1/2} \operatorname{Im} [m_{1c}(w) + m_{2c}(w)] + N^{C_0\delta} \Phi(w). \quad (6.13)$$

$(\mathbf{C}_m)$  For all  $w \in \hat{\mathbf{S}}_m$ , all deterministic unit vector  $\mathbf{v}$ , and all  $X$  satisfying (2.2)-(2.3), we have

$$\left| \mathcal{G}_{\mathbf{v}\mathbf{v}}(w) - \tilde{\Pi}_{\mathbf{v}\mathbf{v}}(w) \right| < N^{C_0\delta} \Phi(w). \quad (6.14)$$

It is trivial to see that property  $(\mathbf{A}_0)$  holds. Moreover, it is easy to observe the following result.

**Lemma 6.4.** For any  $m$ , property  $(\mathbf{C}_m)$  implies property  $(\mathbf{A}_m)$ .

*Proof.* This result follows from (3.33). □

The key step is the following induction result.

**Lemma 6.5.** For any  $1 \leq m \leq 2\delta^{-1}$ , property  $(\mathbf{A}_{m-1})$  implies property  $(\mathbf{C}_m)$ .

Combining Lemmas 6.4 and 6.5, we conclude that (6.14) holds for all  $w \in \hat{\mathbf{S}}$ . Since  $\delta$  can be chosen arbitrarily small under the condition (6.11), we conclude that (6.10) holds for all  $w \in \hat{\mathbf{S}}$ , and Proposition 6.2 follows. What remains now is the proof of Lemma 6.5. Denote

$$F_{\mathbf{v}}(X, w) = \left| \mathcal{G}_{\mathbf{v}\mathbf{v}}(X, w) - \tilde{\Pi}_{\mathbf{v}\mathbf{v}}(w) \right|. \quad (6.15)$$

By Markov's inequality, it suffices to prove the following lemma.

**Lemma 6.6.** Fix  $p \in 2\mathbb{N}$  and  $m \leq 2\delta^{-1}$ . Suppose that the assumptions of Proposition 6.2, (2.23) and property  $(\mathbf{A}_{m-1})$  hold. Then we have

$$\mathbb{E} F_{\mathbf{v}}^p(X, w) \leq \left( N^{C_0\delta} \Phi(w) \right)^p \quad (6.16)$$

for all  $w \in \hat{\mathbf{S}}_m$  and all deterministic unit vector  $\mathbf{v}$ .

In the following, we prove Lemma 6.6. First, in order to make use of the assumption  $(\mathbf{A}_{m-1})$ , which has spectral parameters in  $\hat{\mathbf{S}}_{m-1}$ , to get some estimates for spectral parameters in  $\hat{\mathbf{S}}_m$ , we shall use the following rough bounds for  $\mathcal{G}_{\mathbf{x}\mathbf{y}}$ .

**Lemma 6.7.** For any  $w = E + i\eta \in \mathbf{S}$  and  $\mathbf{x}, \mathbf{y} \in \mathbb{C}^{\mathcal{I}}$ , we have

$$\begin{aligned} \left| \mathcal{G}_{\mathbf{xy}}(w) - \tilde{\Pi}_{\mathbf{xy}}(w) \right| &< N^{2\delta} \sum_{l=1}^{L(\eta)} [\operatorname{Im} \mathcal{G}_{\mathbf{x}_1 \mathbf{x}_1}(E + i\eta_l) + \operatorname{Im} \mathcal{G}_{\mathbf{x}_2 \mathbf{x}_2}(E + i\eta_l) \\ &\quad + \operatorname{Im} \mathcal{G}_{\mathbf{y}_1 \mathbf{y}_1}(E + i\eta_l) + \operatorname{Im} \mathcal{G}_{\mathbf{y}_2 \mathbf{y}_2}(E + i\eta_l)] + |\mathbf{x}| |\mathbf{y}|, \end{aligned}$$

where  $\mathbf{x} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}$  and  $\mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix}$  for  $\mathbf{x}_1, \mathbf{y}_1 \in \mathbb{C}^{\mathcal{I}_1}$  and  $\mathbf{x}_2, \mathbf{y}_2 \in \mathbb{C}^{\mathcal{I}_2}$ .

*Proof.* The proof is similar to the one for [24, Lemma 7.12].  $\square$

**Lemma 6.8.** Suppose  $(\mathbf{A}_{m-1})$  holds, then

$$\mathcal{G}(w) - \tilde{\Pi}(w) = O_{<}(N^{2\delta}) \quad (6.17)$$

and

$$\operatorname{Im} \mathcal{G}_{\mathbf{vv}} \leq N^{2\delta} \left[ |w|^{1/2} \operatorname{Im} (m_{1c}(w) + m_{2c}(w)) + N^{C_0\delta} \Phi(w) \right] \quad (6.18)$$

for all  $w \in \hat{\mathbf{S}}_m$  and all deterministic unit vector  $\mathbf{v}$

*Proof.* Let  $w = E + i\eta \in \hat{\mathbf{S}}_m$ . Then  $E + i\eta_l \in \hat{\mathbf{S}}_{m-1}$  for  $l = 1, \dots, L(\eta)$ , and (6.13) gives  $\operatorname{Im} \mathcal{G}_{\mathbf{vv}}(w) < 1$ . The estimate (6.17) now follows immediately from Lemma 6.7. To prove (6.18), we remark that if  $s(w)$  is the Stieltjes transform of any positive integrable function on  $\mathbb{R}$ , the map  $\eta \mapsto \eta \operatorname{Im} s(E + i\eta)$  is nondecreasing and the map  $\eta \mapsto \eta^{-1} \operatorname{Im} s(E + i\eta)$  is nonincreasing. We apply them to  $|w|^{-1/2} \operatorname{Im} \mathcal{G}_{\mathbf{vv}}(E + i\eta)$  and  $\operatorname{Im} m_{1,2c}(E + i\eta)$  to get for  $w_1 = E + i\eta_1 \in \hat{\mathbf{S}}_{m-1}$ ,

$$\begin{aligned} \operatorname{Im} \mathcal{G}_{\mathbf{vv}}(w) &\leq N^\delta \frac{|w|^{1/2}}{|w_1|^{1/2}} \operatorname{Im} \mathcal{G}_{\mathbf{vv}}(w_1) < N^\delta \left[ |w|^{1/2} \operatorname{Im} (m_{1c}(w_1) + m_{2c}(w_1)) + N^{C_0\delta} \frac{|w|^{1/2}}{|w_1|^{1/2}} \Phi(w_1) \right] \\ &\leq N^{2\delta} \left[ |w|^{1/2} \operatorname{Im} (m_{1c}(w) + m_{2c}(w)) + N^{C_0\delta} \Phi(w) \right], \end{aligned}$$

where we use  $\Phi(w) := |w|^{1/2} \Psi(w)$  and the fact that  $\eta \mapsto \Psi(E + i\eta)$  is nonincreasing, which is clear from the definition (2.45).  $\square$

Now we apply the self-consistent comparison method presented in [24, Section 7] to prove Lemma 6.6. To organize the proof, we divide it into two small subsections.

### 6.1.1 Interpolation and expansion

**Definition 6.9** (Interpolating matrices). Introduce the notation  $X^0 := X^{\text{Gauss}}$  and  $X^1 := X$ . Let  $\rho_{i\mu}^0$  and  $\rho_{i\mu}^1$  be the laws of  $X_{i\mu}^0$  and  $X_{i\mu}^1$ , respectively, for  $i \in \mathcal{I}_1^M$  and  $\mu \in \mathcal{I}_2$ . For  $\theta \in [0, 1]$ , we define the interpolated law

$$\rho_{i\mu}^\theta := (1 - \theta) \rho_{i\mu}^0 + \theta \rho_{i\mu}^1.$$

We shall work on the probability space consisting of triples  $(X^0, X^\theta, X^1)$  of independent  $\mathcal{I}_1^M \times \mathcal{I}_2$  random matrices, where the matrix  $X^\theta = (X_{i\mu}^\theta)$  has law

$$\prod_{i \in \mathcal{I}_1^M} \prod_{\mu \in \mathcal{I}_2} \rho_{i\mu}^\theta(dX_{i\mu}^\theta). \quad (6.19)$$

For  $\lambda \in \mathbb{R}$ ,  $i \in \mathcal{I}_1^M$  and  $\mu \in \mathcal{I}_2$ , we define the matrix  $X_{(i\mu)}^{\theta, \lambda}$  through

$$\left(X_{(i\mu)}^{\theta, \lambda}\right)_{j\nu} := \begin{cases} X_{i\mu}^\theta & \text{if } (j, \nu) \neq (i, \mu) \\ \lambda & \text{if } (j, \nu) = (i, \mu) \end{cases}.$$

We also introduce the matrices

$$\mathcal{G}^\theta(w) := \mathcal{G}(X^\theta, w), \quad \mathcal{G}_{(i\mu)}^{\theta, \lambda}(w) := \mathcal{G}(X_{(i\mu)}^{\theta, \lambda}, w),$$

according to (6.2) and the Definition 2.11.

We shall prove Lemma 6.6 through interpolation matrices  $X^\theta$  between  $X^0$  and  $X^1$ . It holds for  $X^0$  by the anisotropic law (6.3) (see the remark above (6.5)).

**Lemma 6.10.** *Lemma 6.6 holds if  $X = X^0$ .*

Using (6.19) and fundamental calculus, we get the following basic interpolation formula.

**Lemma 6.11.** *For  $F : \mathbb{R}^{\mathcal{I}_1^M \times \mathcal{I}_2} \rightarrow \mathbb{C}$  we have*

$$\frac{d}{d\theta} \mathbb{E}F(X^\theta) = \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \left[ \mathbb{E}F\left(X_{(i\mu)}^{\theta, X_{i\mu}^1}\right) - \mathbb{E}F\left(X_{(i\mu)}^{\theta, X_{i\mu}^0}\right) \right] \quad (6.20)$$

provided all the expectations exists.

We shall apply Lemma 6.11 with  $F(X) = F_{\mathbf{v}}^p(X, w)$  for  $F_{\mathbf{v}}(X, w)$  defined in (6.15). The main work is devoted to prove the following self-consistent estimate for the right-hand side of (6.20).

**Lemma 6.12.** *Fix  $p \in 2\mathbb{N}$  and  $m \leq 2\delta^{-1}$ . Suppose (2.23) and  $(\mathbf{A}_{m-1})$  holds, then we have*

$$\sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \left[ \mathbb{E}F_{\mathbf{v}}^p\left(X_{(i\mu)}^{\theta, X_{i\mu}^1}\right) - \mathbb{E}F_{\mathbf{v}}^p\left(X_{(i\mu)}^{\theta, X_{i\mu}^0}\right) \right] = O\left((N^{C_0\delta}\Phi)^p + \mathbb{E}F_{\mathbf{v}}^p(X^\theta, w)\right) \quad (6.21)$$

for all  $\theta \in [0, 1]$ , all  $w \in \hat{\mathbf{S}}_m$ , and all deterministic unit vector  $\mathbf{v}$ .

Combining Lemmas 6.10, 6.11 and 6.12 with a Grönwall argument, we can conclude the proof of Lemma 6.6 and hence Proposition 6.2.

In order to prove Lemma 6.12, we compare  $X_{(i\mu)}^{\theta, X_{i\mu}^0}$  and  $X_{(i\mu)}^{\theta, X_{i\mu}^1}$  via a common  $X_{(i\mu)}^{\theta, 0}$ , i.e. under the assumptions of Lemma 6.12, we will prove

$$\sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \left[ \mathbb{E}F_{\mathbf{v}}^p\left(X_{(i\mu)}^{\theta, X_{i\mu}^u}\right) - \mathbb{E}F_{\mathbf{v}}^p\left(X_{(i\mu)}^{\theta, 0}\right) \right] = O\left((N^{C_0\delta}\Phi)^p + \mathbb{E}F_{\mathbf{v}}^p(X^\theta, w)\right) \quad (6.22)$$

for all  $u \in \{0, 1\}$ , all  $\theta \in [0, 1]$ , all  $w \in \hat{\mathbf{S}}_m$ , and all deterministic unit vector  $\mathbf{v}$ .

Underlying the proof of (6.22) is an expansion approach which we will describe below. Throughout the rest of the proof, we suppose that  $(\mathbf{A}_{m-1})$  holds. Also the rest of the proof is performed at a single  $w \in \hat{\mathbf{S}}_m$ . Define the  $\mathcal{I} \times \mathcal{I}$  matrix  $\Delta_{(i\mu)}^\lambda$  through

$$\left(\Delta_{(i\mu)}^\lambda\right)_{st} := \lambda \delta_{is} \delta_{\mu t} + \lambda \delta_{it} \delta_{\mu s}. \quad (6.23)$$

Then we have for any  $\lambda, \lambda' \in \mathbb{R}$  and  $K \in \mathbb{N}$ ,

$$\mathcal{G}_{(i\mu)}^{\theta, \lambda'} = \mathcal{G}_{(i\mu)}^{\theta, \lambda} + \sum_{k=1}^K \alpha^k \mathcal{G}_{(i\mu)}^{\theta, \lambda} \left( \bar{V} \Delta_{(i\mu)}^{\lambda - \lambda'} \bar{V}^\dagger \mathcal{G}_{(i\mu)}^{\theta, \lambda} \right)^k + \alpha^{K+1} \mathcal{G}_{(i\mu)}^{\theta, \lambda'} \left( \bar{V} \Delta_{(i\mu)}^{\lambda - \lambda'} \bar{V}^\dagger \mathcal{G}_{(i\mu)}^{\theta, \lambda} \right)^{K+1}, \quad (6.24)$$

where  $\bar{V} := \begin{pmatrix} V_1 & 0 \\ 0 & I \end{pmatrix}$  and  $\alpha := \frac{w^{1/2}}{|w|^{1/2}}$ . The following result provides a priori bounds for the entries of  $\mathcal{G}_{(i\mu)}^{\theta, \lambda}$ .

**Lemma 6.13.** *Suppose that  $y$  is a random variable satisfying  $|y| < N^{-1/2}$ . Then*

$$\mathcal{G}_{(i\mu)}^{\theta, y} - \tilde{\Pi} = O_{<}(N^{2\delta}) \quad (6.25)$$

for all  $i \in \mathcal{I}_1^M$  and  $\mu \in \mathcal{I}_2$ .

*Proof.* See [24, Lemma 7.14]. □

In the following, for simplicity of notations we introduce  $f_{(i\mu)}(\lambda) := F_{\mathbf{v}}^p(X_{(i\mu)}^{\theta, \lambda})$ . We use  $f_{(i\mu)}^{(n)}$  to denote the  $n$ -th derivative of  $f_{(i\mu)}$ . By Lemma 6.13 and expansion (6.24) we get the following result.

**Lemma 6.14.** *Suppose that  $y$  is a random variable satisfying  $|y| < N^{-1/2}$ . Then for fixed  $n \in \mathbb{N}$ ,*

$$\left| f_{(i\mu)}^{(n)}(y) \right| < N^{2\delta(n+p)}. \quad (6.26)$$

By this lemma, the Taylor expansion of  $f_{(i\mu)}$  gives

$$f_{(i\mu)}(y) = \sum_{n=0}^{4p} \frac{y^n}{n!} f_{(i\mu)}^{(n)}(0) + O_{<}(\Phi^p), \quad (6.27)$$

provided  $C_0$  is chosen large enough in (6.11). Therefore we have for  $u \in \{0, 1\}$ ,

$$\begin{aligned} \mathbb{E} F_{\mathbf{v}}^p \left( X_{(i\mu)}^{\theta, X_{i\mu}^u} \right) - \mathbb{E} F_{\mathbf{v}}^p \left( X_{(i\mu)}^{\theta, 0} \right) &= \mathbb{E} \left[ f_{(i\mu)} \left( X_{i\mu}^u \right) - f_{(i\mu)}(0) \right] \\ &= \mathbb{E} f_{(i\mu)}(0) + \frac{1}{2N} \mathbb{E} f_{(i\mu)}^{(2)}(0) + \sum_{n=4}^{4p} \frac{1}{n!} \mathbb{E} f_{(i\mu)}^{(n)}(0) \mathbb{E} \left( X_{i\mu}^u \right)^n + O_{<}(\Phi^p), \end{aligned}$$

where we used that  $X_{i\mu}^u$  has vanishing first and third moments and its variance is  $1/N$ . Thus to show (6.22), we only need to prove for  $n = 4, 5, \dots, 4p$ ,

$$N^{-n/2} \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \left| \mathbb{E} f_{(i\mu)}^{(n)}(0) \right| = O \left( (N^{C_0 \delta} \Phi)^p + \mathbb{E} F_{\mathbf{v}}^p(X^\theta, w) \right), \quad (6.28)$$

where we have used (2.3). In order to get a self-consistent estimate in terms of the matrix  $X^\theta$  on the right-hand side of (6.28), we want to replace  $X_{(i\mu)}^{\theta, 0}$  in  $f_{(i\mu)}(0) := F_{\mathbf{v}}^p(X_{(i\mu)}^{\theta, 0})$  with  $X^\theta = X_{(i\mu)}^{\theta, X_{i\mu}^\theta}$ .

**Lemma 6.15.** *Suppose that*

$$N^{-n/2} \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \left| \mathbb{E} f_{(i\mu)}^{(n)}(X_{i\mu}^\theta) \right| = O \left( (N^{C_0 \delta} \Phi)^p + \mathbb{E} F_{\mathbf{v}}^p(X^\theta, w) \right) \quad (6.29)$$

holds for  $n = 4, \dots, 4p$ , Then (6.28) holds for  $n = 4, \dots, 4p$ .

*Proof.* From (6.27) we can get

$$f_{(i\mu)}^{(l)}(0) = f_{(i\mu)}^{(l)}(y) - \sum_{n=1}^{4p-l} \frac{y^n}{n!} f_{(i\mu)}^{(l+n)}(0) + O_{<}(N^{l/2}\Phi^p). \quad (6.30)$$

The result follows by repeatedly applying (6.30). The details can be found in [24, Lemma 7.16].  $\square$

### 6.1.2 Conclusion of the proof with words

What remains now is to prove (6.29). In order to exploit the detailed structure of the derivatives on the left-hand side of (6.29), we introduce the following algebraic objects.

**Definition 6.16** (Words). *Given  $i \in \mathcal{I}_1^M$  and  $\mu \in \mathcal{I}_2$ . Let  $\mathcal{W}$  be the set of words of even length in two letters  $\{\mathbf{i}, \mu\}$ . We denote the length of a word  $w \in \mathcal{W}$  by  $2n(w)$  with  $n(w) \in \mathbb{N}$ . We use bold symbols to denote the letters of words. For instance,  $w = \mathbf{t}_1 \mathbf{s}_2 \mathbf{t}_2 \mathbf{s}_3 \cdots \mathbf{t}_n \mathbf{s}_{n+1}$  denotes a word of length  $2n$ . Define  $\mathcal{W}_n := \{w \in \mathcal{W} : n(w) = n\}$  to be the set of words of length  $2n$ . We require that each word  $w \in \mathcal{W}_n$  satisfies that  $\mathbf{t}_l \mathbf{s}_{l+1} \in \{\mathbf{i}\mu, \mu\mathbf{i}\}$  for all  $1 \leq l \leq n$ .*

*Next we assign each letter  $*$  its value  $[*]$  through  $[\mathbf{i}] := \mathbf{v}_i$ ,  $[\mu] := \mu$ , where  $\mathbf{v}_i \in \mathbb{C}^{\mathcal{I}_1}$  is defined in Lemma 6.1 and is regarded as a summation index. Note that it is important to distinguish the abstract letter from its value, which is a summation index. Finally, to each word  $w$  we assign a random variable  $A_{\mathbf{v}, i, \mu}(w)$  as follows. If  $n(w) = 0$  we define*

$$A_{\mathbf{v}, i, \mu}(W) := \mathcal{G}_{\mathbf{v}\mathbf{v}} - \tilde{\Pi}_{\mathbf{v}\mathbf{v}}.$$

*If  $n(w) \geq 1$ , say  $w = \mathbf{t}_1 \mathbf{s}_2 \mathbf{t}_2 \mathbf{s}_3 \cdots \mathbf{t}_n \mathbf{s}_{n+1}$ , we define*

$$A_{\mathbf{v}, i, \mu}(W) := \mathcal{G}_{\mathbf{v}[\mathbf{t}_1]} \mathcal{G}_{[\mathbf{s}_2][\mathbf{t}_2]} \cdots \mathcal{G}_{[\mathbf{s}_n][\mathbf{t}_n]} \mathcal{G}_{[\mathbf{s}_{n+1}]} \mathbf{v}. \quad (6.31)$$

Notice the words are constructed such that, by (6.24),

$$\left( \frac{\partial}{\partial X_{i\mu}} \right)^n \left( \mathcal{G}_{\mathbf{v}\mathbf{v}} - \tilde{\Pi}_{\mathbf{v}\mathbf{v}} \right) = (-\alpha)^n n! \sum_{w \in \mathcal{W}_n} A_{\mathbf{v}, i, \mu}(w)$$

for  $n = 0, 1, 2, \dots$ , which gives that

$$\begin{aligned} \left( \frac{\partial}{\partial X_{i\mu}} \right)^n F_{\mathbf{v}}^p(X) &= (-\alpha)^n n! \sum_{n_1 + \cdots + n_p = n} \prod_{r=1}^{p/2} \frac{1}{n_r! n_{r+p/2}!} \\ &\quad \times \left( \sum_{w_r \in \mathcal{W}_{n_r}} \sum_{w_{r+p/2} \in \mathcal{W}_{n_{r+p/2}}} A_{\mathbf{v}, i, \mu}(w_r) \overline{A_{\mathbf{v}, i, \mu}(w_{r+p/2})} \right). \end{aligned}$$

Then to prove (6.29), it suffices to show that

$$N^{-n/2} \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \left| \mathbb{E} \prod_{r=1}^{p/2} A_{\mathbf{v}, i, \mu}(w_r) \overline{A_{\mathbf{v}, i, \mu}(w_{r+p/2})} \right| = O \left( (N^{C_0 \delta} \Phi)^p + \mathbb{E} F_{\mathbf{v}}^p(X^\theta, w) \right) \quad (6.32)$$

for  $4 \leq n \leq 4p$  and all words  $w_1, \dots, w_p \in \mathcal{W}$  satisfying  $n(w_1) + \cdots + n(w_p) = n$ . To avoid the unimportant notational complications coming from the complex conjugates, we in fact prove that

$$N^{-n/2} \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \left| \mathbb{E} \prod_{r=1}^p A_{\mathbf{v}, i, \mu}(w_r) \right| = O \left( (N^{C_0 \delta} \Phi)^p + \mathbb{E} F_{\mathbf{v}}^p(X^\theta, w) \right), \quad (6.33)$$

and the proof of (6.32) is essentially the same but with slightly heavier notations. Treating empty words separately, we find it suffices to prove

$$N^{-n/2} \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \mathbb{E} \left| A_{\mathbf{v}, i, \mu}^{p-q}(w_0) \prod_{r=1}^q A_{\mathbf{v}, i, \mu}(w_r) \right| = O((N^{C_0 \delta} \Phi)^p + \mathbb{E} F_{\mathbf{v}}^p(X^\theta, w)) \quad (6.34)$$

for  $4 \leq n \leq 4p$ ,  $1 \leq q \leq p$ , and  $w_r$  such that  $n(w_0) = 0$ ,  $\sum_r n(w_r) = n$  and  $n(w_r) \geq 1$  for  $r \geq 1$ .

To estimate (6.34) we introduce the quantity

$$\mathcal{R}_s := |\mathcal{G}_{\mathbf{v}\mathbf{v}_s}| + |\mathcal{G}_{\mathbf{v}_s\mathbf{v}}|. \quad (6.35)$$

for  $s \in \mathcal{I}$ , where as a convention we let  $\mathbf{v}_\mu = e_\mu$  for  $\mu \in \mathcal{I}_2$ .

**Lemma 6.17.** *For  $w \in \mathcal{W}$  we have the rough bound*

$$|A_{\mathbf{v}, i, \mu}(w)| < N^{2\delta(n(w)+1)}. \quad (6.36)$$

Furthermore, for  $n(w) \geq 1$  we have

$$|A_{\mathbf{v}, i, \mu}(w)| < (\mathcal{R}_i^2 + \mathcal{R}_\mu^2) N^{2\delta(n(w)-1)}. \quad (6.37)$$

For  $n(w) = 1$  we have better bound

$$|A_{\mathbf{v}, i, \mu}(w)| < \mathcal{R}_i \mathcal{R}_\mu. \quad (6.38)$$

*Proof.* (6.36) follows immediately from the rough bound (6.17) and definition (6.31). For (6.37) we break  $A_{\mathbf{v}, i, \mu}(w)$  into  $\mathcal{G}_{\mathbf{v}[\mathbf{t}_1]}(\mathcal{G}_{[\mathbf{s}_2][\mathbf{t}_2]} \cdots \mathcal{G}_{[\mathbf{s}_n][\mathbf{t}_n]})^{1/2}$  times  $(\mathcal{G}_{[\mathbf{s}_2][\mathbf{t}_2]} \cdots \mathcal{G}_{[\mathbf{s}_n][\mathbf{t}_n]})^{1/2} \mathcal{G}_{[\mathbf{s}_{n+1}]\mathbf{v}}$  and use Cauchy-Schwarz inequality. (6.38) follows from the constraint  $\mathbf{t}_1 \neq \mathbf{s}_2$  in the definition (6.31).  $\square$

By pigeonhole principle, if  $n \leq 2q - 2$  there exists at least two words  $w_r$  with  $n(w_r) = 1$ . Therefore by Lemma 6.17 we have

$$\left| A_{\mathbf{v}, i, \mu}^{p-q}(w_0) \prod_{r=1}^q A_{\mathbf{v}, i, \mu}(w_r) \right| < N^{2\delta(n+q)} F_{\mathbf{v}}^{p-q}(X) \left( \mathbf{1}(n \geq 2q - 1) (\mathcal{R}_i^2 + \mathcal{R}_\mu^2) + \mathbf{1}(n \leq 2q - 2) \mathcal{R}_i^2 \mathcal{R}_\mu^2 \right). \quad (6.39)$$

Then by Lemma 6.1,

$$\begin{aligned} \frac{1}{N} \sum_{i \in \mathcal{I}_1^M} \mathcal{R}_i^2 + \frac{1}{N} \sum_{\mu \in \mathcal{I}_2} \mathcal{R}_\mu^2 &< \frac{|w|^{1/2} \text{Im} \mathcal{G}_{\mathbf{v}\mathbf{v}} + \eta |w|^{-1/2} \mathcal{G}_{\mathbf{v}\mathbf{v}}}{N\eta} \\ &< N^{2\delta} \frac{|w| \text{Im}(m_{1c} + m_{2c}) + |w|^{1/2} N^{C_0 \delta} \Phi}{N\eta} < N^{(C_0+2)\delta} \Phi^2, \end{aligned} \quad (6.40)$$

where in the second step we used the two bounds in Lemma 6.8,  $|w|^{-1/2} \eta = O(|w| \text{Im} m_{1c})$  by Lemma 3.7, and in the last step the definition of  $\Phi$ . Using the same method we can get

$$\frac{1}{N^2} \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \mathcal{R}_i^2 \mathcal{R}_\mu^2 < \left( N^{(C_0+2)\delta} \Phi^2 \right)^2. \quad (6.41)$$

Plugging (6.40) and (6.41) into (6.39), we get that the left-hand side of (6.34) is bounded by

$$N^{-n/2+2} N^{2\delta(n+q+2)} \mathbb{E} F_{\mathbf{v}}^{p-q}(X) \left( \mathbf{1}(n \geq 2q - 1) \left( N^{C_0 \delta/2} \Phi \right)^2 + \mathbf{1}(n \leq 2q - 2) \left( N^{C_0 \delta/2} \Phi \right)^4 \right).$$

Using  $\Phi \geq cN^{-1/2}$ , we find that the left hand side of (6.34) is bounded by

$$\begin{aligned} & N^{2\delta(n+q+2)} \mathbb{E} F_{\mathbf{v}}^{p-q}(X) \left( \mathbf{1}(n \geq 2q-1) \left( N^{C_0\delta/2} \Phi \right)^{n-2} + \mathbf{1}(n \leq 2q-2) \left( N^{C_0\delta/2} \Phi \right)^n \right) \\ & \leq \mathbb{E} F_{\mathbf{v}}^{p-q}(X) \left( \mathbf{1}(n \geq 2q-1) \left( N^{C_0\delta/2+12\delta} \Phi \right)^{n-2} + \mathbf{1}(n \leq 2q-2) \left( N^{C_0\delta/2+12\delta} \Phi \right)^n \right) \end{aligned}$$

where we used that  $q \leq n$  and  $n \geq 4$ . Choose  $C_0 \geq 25$ , then by (6.11) we have  $N^{C_0\delta/2+12\delta} \leq N^{\zeta/2}$  and hence  $N^{C_0\delta/2+12\delta} \Phi \leq 1$ . Moreover, if  $n \geq 4$  and  $n \geq 2q-1$ , then  $n \geq q+2$ . Therefore we conclude that the left-hand side of (6.34) is bounded by

$$\mathbb{E} F_{\mathbf{v}}^{p-q}(X) (N^{C_0\delta} \Phi)^q. \quad (6.42)$$

Now (6.34) follows from Holder's inequality. This concludes the proof of (6.29), and hence of (6.22), and then of Lemma 6.5. This finishes the proof of Proposition 6.2 under the assumption (2.23).

In the rest of this section, we prove Proposition 6.2 when  $\eta \geq N^{-1/2+\zeta} |m_{2c}|^{-1}$ . In this case, we can verify that

$$\Phi \leq N^{-1/4-\zeta/2}. \quad (6.43)$$

Following the previous arguments, we see that it suffices to prove the estimate (6.29) for  $n = 3$ . In other words, we need to prove the following lemma.

**Lemma 6.18.** *Fix  $1 \leq m \leq 2\delta^{-1}$  and  $p \in 2\mathbb{N}$ . Let  $w \in \hat{\mathbf{S}}_m \cap \hat{\mathbf{D}}$  (recall (2.44)) and suppose  $(\mathbf{A}_{m-1})$  holds. Then we have*

$$N^{-3/2} \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \left| \mathbb{E} f_{(i\mu)}^{(3)}(X_{i\mu}^\theta) \right| = O \left( (N^{C_0\delta} \Phi)^p + \mathbb{E} F_{\mathbf{v}}^p(X^\theta, w) \right). \quad (6.44)$$

*Proof.* The main new ingredient of the proof is a further iteration step at a fixed  $w$ . Suppose

$$\mathcal{G} - \tilde{\Pi} = O_{<}(N^{2\delta} \phi) \quad (6.45)$$

for some  $\phi \leq 1$ . By the a priori bound (6.17), (6.45) holds for  $\phi = 1$ . Assuming (6.45), we shall prove a self-improving bound of the form

$$N^{-3/2} \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \left| \mathbb{E} f_{(i\mu)}^{(3)}(X_{i\mu}^\theta) \right| = O \left( (N^{C_0\delta} \Phi)^p + (N^{-\zeta/4} \phi)^p + \mathbb{E} F_{\mathbf{v}}^p(X^\theta, w) \right). \quad (6.46)$$

Once (6.46) is proved, we can use it iteratively to get an increasingly accurate bound for the left hand side of (6.14). After each step, we obtain a better a priori bound (6.45) where  $\phi$  is reduced by  $N^{-\zeta/4}$ . Hence after  $O(\zeta^{-1})$  iterations we can get (6.44).

As in Section 6.1.2, to prove (6.46) it suffice to show

$$N^{-3/2} \left| \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} A_{\mathbf{v}, i, \mu}^{p-q}(w_0) \prod_{r=1}^q A_{\mathbf{v}, i, \mu}(w_r) \right| < F_{\mathbf{v}}^{p-q}(X) (N^{(C_0-1)\delta} \Phi + N^{-\zeta/2} \phi)^q, \quad (6.47)$$

which follows from

$$N^{-3/2} \left| \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \prod_{r=1}^q A_{\mathbf{v}, i, \mu}(w_r) \right| < (N^{(C_0-1)\delta} \Phi + N^{-\zeta/2} \phi)^q. \quad (6.48)$$

Each of the three cases  $q = 1, 2, 3$  can be proved as in [24, Lemma 12.7], and we leave the details to the reader. This concludes Lemma 6.18.  $\square$

## 6.2 Averaged local law for $TX$

In this section we prove the averaged local law in Theorem 2.19. Again for convenience, we only consider the case  $w \in \mathbf{D}$  and  $|z|^2 \leq 1 - \tau$ . First we assume (2.23) holds. The anisotropic local law proved in the previous section gives a good a priori bound. In analogy to (6.15), we define

$$\tilde{F}(X, w) := |w|^{1/2} |m_2(w) - m_{2c}(w)| = \left| \frac{1}{N} \sum_{\nu \in \mathcal{I}_2} \mathcal{G}_{\nu\nu}(w) - |w|^{1/2} m_{2c}(w) \right|.$$

Since  $\Phi^2 = O(|w|^{1/2}/(N\eta))$ , it suffices to show that  $\tilde{F} < \Phi^2$ . Following the argument in Section 6.1, analogous to (6.29), we only need to prove that

$$N^{-n/2} \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \left| \mathbb{E} \left( \frac{\partial}{\partial X_{i\mu}} \right)^n \tilde{F}^p(X) \right| = O \left( (N^\delta \Phi^2)^p + \mathbb{E} \tilde{F}^p(X) \right) \quad (6.49)$$

for all  $n = 4, \dots, 4p$ . Here  $\delta > 0$  is an arbitrary positive constant. Analogously to (6.33), it suffices to prove that for  $n = 4, \dots, 4p$ ,

$$N^{-n/2} \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \left| \mathbb{E} \prod_{r=1}^p \left( \frac{1}{N} \sum_{\nu \in \mathcal{I}_2} A_{\mathbf{e}_\nu, i, \mu}(w_r) \right) \right| = O \left( (N^\delta \Phi^2)^p + \mathbb{E} \tilde{F}^p(X) \right) \quad (6.50)$$

for  $\sum_r n(w_r) = n$ . The only difference in the definition of  $A_{\mathbf{v}, i, \mu}(w)$  is that when  $n(w) = 0$ , we define

$$A_{\mathbf{v}, i, \mu}(w) := \mathcal{G}_{\mathbf{v}\mathbf{v}} - |w|^{1/2} m_{2c}.$$

Similar to (6.35) we define

$$\mathcal{R}_{\nu, s} := |\mathcal{G}_{\nu\mathbf{v}_s}| + |\mathcal{G}_{\mathbf{v}_s\nu}|. \quad (6.51)$$

By the anisotropic local law,  $\mathcal{G} - \tilde{\Pi} = O_{<}(\Phi)$ . Hence combining with Lemma 6.1 and (3.33), we get

$$\frac{1}{N} \sum_{\nu \in \mathcal{I}_2} \mathcal{R}_{\nu, s}^2 < \frac{|w|^{1/2} \text{Im} \mathcal{G}_{\mathbf{v}_s \mathbf{v}_s}}{N\eta} < \frac{|w| \text{Im}(m_{1c} + m_{2c}) + |w|^{1/2} \Phi}{N\eta} = O(\Phi^2). \quad (6.52)$$

Using the anisotropic local law again, we get  $\mathcal{G} = O_{<}(1)$ . Then we have

$$\left| \frac{1}{N} \sum_{\nu \in \mathcal{I}_2} A_{\mathbf{e}_\nu, i, \mu}(w) \right| < \frac{1}{N} \sum_{\nu \in \mathcal{I}_2} (\mathcal{R}_{\nu, i}^2 + \mathcal{R}_{\nu, \mu}^2) < \Phi^2 \text{ for } n(w) \geq 1. \quad (6.53)$$

Following (6.53), for  $n \geq 4$ , the left-hand side of (6.50) is bounded by

$$\mathbb{E} \tilde{F}^{p-q}(X) (\Phi^2)^q.$$

Applying Holder's inequality, we conclude the proof.

Then we prove the averaged local law when  $\eta \geq N^{-1/2+\zeta} |m_{2c}|^{-1}$ . It suffices to prove

$$N^{-3/2} \left| \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \mathbb{E} \left( \frac{\partial}{\partial X_{i\mu}} \right)^3 \tilde{F}^p(X) \right| = O \left( \left( \frac{|w|^{1/2}}{N\eta} \right)^p + \mathbb{E} \tilde{F}^p(X) \right). \quad (6.54)$$



Analogous to (6.50), it is reduced to show that

$$N^{-3/2} \left| \sum_{i \in \mathcal{I}_1^M} \sum_{\mu \in \mathcal{I}_2} \mathbb{E} \prod_{r=1}^q \left( \frac{1}{N} \sum_{\nu \in \mathcal{I}_2} A_{\mathbf{e}_\nu, i, \mu}(w_r) \right) \right| = O \left( \left( \frac{|w|^{1/2}}{N\eta} \right)^q + \mathbb{E} \tilde{F}^q(X) \right) \quad (6.55)$$

where  $q$  is the number of words with nonzero length. Again we can prove the three cases  $q = 1, 2, 3$  as in [24, Lemma 12.8], and we leave the details to the reader. This concludes the averaged law.

## A Properties of $\rho_{1,2c}$ and Stability of (2.11)

### A.1 Proof of Lemma 2.3 and Proposition 2.14

We now prove Lemma 2.3. First is a technical lemma for  $f$  defined in (2.15).

**Lemma A.1.** *For  $w > 0$  and  $|z| > 0$ ,  $f$  can be written as*

$$f(\sqrt{w}, m) = -\sqrt{w} + m + w^{-1/2} + \frac{1}{N} \sum_{i=1}^n l_i s_i \left( \frac{A_i}{m - a_i} + \frac{B_i}{m - b_i} + \frac{C_i}{m + c_i} \right), \quad (A.1)$$

where we have the following estimates for the poles and the coefficients,

$$\max \left( |z|, \frac{s_i + |z|^2}{\sqrt{w}} \right) < a_i < \frac{s_i + |z|^2}{\sqrt{w}} + |z|, \quad a_n < a_{n-1} < \dots < a_1, \quad (A.2)$$

$$0 < b_1 < b_2 < \dots < b_n < \min \left( |z|, \frac{|z|^2}{\sqrt{w}} \right), \quad (A.3)$$

$$\frac{-(s_i + |z|^2) + \sqrt{(s_i + |z|^2)^2 + 4w|z|^2}}{2\sqrt{w}} < c_i < |z|, \quad c_1 < c_2 < \dots < c_n, \quad (A.4)$$

and

$$0 < A_i \leq 2 \frac{s_i + |z|^2 + \sqrt{w}|z|}{w}, \quad 0 < B_i \leq 2 \frac{s_i + |z|^2 + \sqrt{w}|z|}{w}, \quad 0 < C_i \leq \frac{s_i + |z|^2 + \sqrt{w}|z|}{w}. \quad (A.5)$$

*Proof.* The proof is based on basic algebraic arguments. Let

$$p_i = \sqrt{w}m^3 - (s_i + |z|^2)m^2 - \sqrt{w}|z|^2m + |z|^4.$$

It is easy to verify that

$$\Delta = 18(s_i + |z|^2)w|z|^6 + 4(s_i + |z|^2)^3|z|^4 + (s_i + |z|^2)^2w|z|^4 + 4w^2|z|^6 - 27w|z|^8 > 0.$$

Thus  $p_i$  has three distinct real roots. By the form of  $p_i$ , we see that there are two positive roots and one negative root, call them  $a_i > b_i > 0 > -c_i$ . Now we perform the partial fraction expansion for the rational functions in (2.15),

$$\frac{m^2 - |z|^2}{\sqrt{w}m^3 - (s_i + |z|^2)m^2 - \sqrt{w}|z|^2m + |z|^4} = \frac{A'_i}{m - a_i} + \frac{B'_i}{m - b_i} - \frac{C'_i}{m + c_i}, \quad (A.6)$$

where

$$A'_i = \frac{a_i^2 - |z|^2}{\sqrt{w}(a_i - b_i)(a_i + c_i)}, \quad B'_i = \frac{b_i^2 - |z|^2}{\sqrt{w}(b_i - a_i)(b_i + c_i)}, \quad C'_i = \frac{-c_i^2 + |z|^2}{\sqrt{w}(c_i + a_i)(c_i + b_i)}. \quad (A.7)$$

We take  $s_i = 0$  in  $p_i$  and call the resulting polynomial as

$$p_0 = \sqrt{w}m^3 - |z|^2m^2 - \sqrt{w}|z|^2m + |z|^4 = \sqrt{w} \left( m - \frac{|z|^2}{\sqrt{w}} \right) (m^2 - |z|^2),$$

which has roots  $m = \pm|z|, |z|^2/\sqrt{w}$ . By (2.7), we have  $p_1 < p_2 < \dots < p_n < p_0$  for all  $m \neq 0$ . Comparing the graphs of  $p_i$ 's (as cubic functions of  $m$ ) for  $0 \leq i \leq n$ , we get that

$$\max \left( |z|, \frac{|z|^2}{\sqrt{w}} \right) < a_n < a_{n-1} < \dots < a_1, \quad 0 < b_1 < b_2 < \dots < b_n < \min \left( |z|, \frac{|z|^2}{\sqrt{w}} \right), \quad (\text{A.8})$$

and

$$0 < c_1 < c_2 < \dots < c_n < |z|. \quad (\text{A.9})$$

Thus we get (A.3). By these bounds, we see that  $a_i^2 - |z|^2 > 0$ ,  $b_i^2 - |z|^2 < 0$  and  $-c_i^2 + |z|^2 > 0$ , which, by (A.7), give that  $A'_i > 0$ ,  $B'_i > 0$  and  $C'_i > 0$ . Plugging (A.6) into  $f$ , we get immediately (A.1) for  $A_i = A'_i a_i$ ,  $B_i = B'_i b_i$  and  $C_i = C'_i c_i$ .

Now we compare  $p_i$  with  $p'_i := \sqrt{w}m^3 - (s_i + |z|^2)m^2 - \sqrt{w}|z|^2m$ , which has roots

$$m = 0, \quad \frac{(s_i + |z|^2) \pm \sqrt{(s_i + |z|^2)^2 + 4w|z|^2}}{2\sqrt{w}}.$$

Since  $p'_i < p_i$  for all  $m$ , we get

$$a_i < \frac{(s_i + |z|^2) + \sqrt{(s_i + |z|^2)^2 + 4w|z|^2}}{2\sqrt{w}} < \frac{s_i + |z|^2}{\sqrt{w}} + |z|, \quad (\text{A.10})$$

and

$$c_i > \frac{-(s_i + |z|^2) + \sqrt{(s_i + |z|^2)^2 + 4w|z|^2}}{2\sqrt{w}}. \quad (\text{A.11})$$

From (A.9) and (A.11), we get (A.4). Then we compare  $p_i$  with  $p''_i := \sqrt{w}m^3 - (s_i + |z|^2)m^2$ , which has roots  $w = 0, (s_i + |z|^2)/\sqrt{w}$ . Notice  $p''_i > p_i$  for  $m > |z|^2/\sqrt{w}$  and  $a_i > |z|^2/\sqrt{w}$ , so we get  $a_i > (s_i + |z|^2)/\sqrt{w}$ . Combining this bound with (A.8) and (A.10), we get (A.2).

Finally we estimate the coefficients  $A_i$ ,  $B_i$  and  $C_i$ . Using (A.7) and (A.2)-(A.4), we first can estimate that

$$\begin{aligned} A'_i &= \frac{(a_i - |z|)(a_i + |z|)}{\sqrt{w}(a_i - b_i)(a_i + c_i)} \leq \frac{a_i + |z|}{\sqrt{w}(a_i + c_i)} \leq \frac{2}{\sqrt{w}}, \\ B'_i &= \frac{(|z| + b_i)(|z| - b_i)}{\sqrt{w}(a_i - b_i)(b_i + c_i)} \leq \frac{|z| + b_i}{\sqrt{w}(b_i + c_i)} \leq 2 \frac{s_i + |z|^2 + \sqrt{w}|z|}{w|z|}, \\ C'_i &= \frac{(|z| - c_i)(c_i + |z|)}{\sqrt{w}(c_i + a_i)(c_i + b_i)} \leq \frac{|z| - c_i}{\sqrt{w}(c_i + b_i)} \leq \frac{s_i + |z|^2 + \sqrt{w}|z|}{w|z|}, \end{aligned}$$

from which we get that

$$A_i = A'_i a_i \leq \frac{2}{\sqrt{w}} \left( \frac{s_i + |z|^2}{\sqrt{w}} + |z| \right) = 2 \frac{s_i + |z|^2 + \sqrt{w}|z|}{w}, \quad (\text{A.12})$$

$$B_i = B'_i b_i \leq 2 \frac{s_i + |z|^2 + \sqrt{w}|z|}{w|z|} |z| = 2 \frac{s_i + |z|^2 + \sqrt{w}|z|}{w}, \quad (\text{A.13})$$

$$C_i = C'_i c_i \leq \frac{s_i + |z|^2 + \sqrt{w}|z|}{w|z|} |z| = \frac{s_i + |z|^2 + \sqrt{w}|z|}{w}. \quad (\text{A.14})$$

□

In (A.1), it is sometimes convenient to reorder the terms and rename the constants to write  $f$  as

$$f(m) = -\sqrt{w} + m + w^{-1/2} + \frac{1}{N} \sum_{k=1}^{2n} \frac{C_k^+}{m - x_k} + \frac{1}{N} \sum_{l=1}^n \frac{C_l^-}{m + y_l}. \quad (\text{A.15})$$

where all the constants  $C_k^+$  and  $C_l^-$  are positive, and we choose the order such that

$$0 < x_1 < x_2 < \dots < x_{2n}, \quad 0 < y_1 < y_2 < \dots < y_n. \quad (\text{A.16})$$

Clearly,  $f$  is smooth on the  $3n + 1$  open intervals of  $\mathbb{R}$  defined by

$$I_{-n} := (-\infty, -y_n), \quad I_{-k} := (-y_{k+1}, -y_k) \quad (k = 1, \dots, n-1), \quad I_0 := (-y_1, x_1), \\ I_k := (x_k, x_{k+1}) \quad (k = 1, \dots, 2n-1), \quad I_{2n} := (x_{2n}, +\infty).$$

Next, we introduce the multiset  $\mathcal{C}$  of critical points of  $f$  (as a function of  $m$ ), using the conventions that a nondegenerate critical point is counted once and a degenerated critical point twice. First we will prove the following elementary lemma about the structure of  $\mathcal{C}$  (see Fig. 6 and 7).

**Lemma A.2.** (*Critical points*) *We have  $|\mathcal{C} \cap I_{-n}| = |\mathcal{C} \cap I_{2n}| = 1$  and  $|\mathcal{C} \cap I_k| \in \{0, 2\}$  for  $k = -n + 1, \dots, 2n - 1$ .*

*Proof.* We omit the dependence of  $f$  on  $w$  for now. By (A.15) we have

$$f'(m) = 1 - \frac{1}{N} \sum_{k=1}^{2n} \frac{C_k^+}{(m - x_k)^2} - \frac{1}{N} \sum_{l=1}^n \frac{C_l^-}{(m + y_l)^2}, \quad f''(m) = \frac{1}{N} \sum_{k=1}^{2n} \frac{2C_k^+}{(m - x_k)^3} + \frac{1}{N} \sum_{l=1}^n \frac{2C_l^-}{(m + y_l)^3}.$$

We see that  $f''$  is decreasing on all the intervals  $I_k$  for  $k = -n + 1, \dots, 2n - 1$ . Thus there is at most one point  $m \in I_k$  such that  $f''(m) = 0$ . We conclude that  $f$  has at most two critical points on  $I_k$ . By the boundary conditions of  $f'$  on  $\partial I_k$ , we get  $|\mathcal{C} \cap I_k| \in \{0, 2\}$  for  $k = -n + 1, \dots, 2n - 1$ . For  $m < -y_n$ , we have  $f''(m) < 0$ , while for  $m > x_{2n}$ , we have  $f''(m) > 0$ . By the boundary conditions of  $f'$  on  $\partial I_{-n}$  and  $\partial I_{2n}$ , we see that  $f'$  decreases from 1 to  $-\infty$  when  $m$  increases from  $-\infty$  to  $-y_n$ , while  $f'$  increases from  $-\infty$  to 1 when  $m$  increases from  $x_{2n}$  to  $+\infty$ . Hence we conclude that each of the intervals  $(-\infty, -y_n)$  and  $(x_{2n}, +\infty)$  contains a unique critical point in it, i.e.  $|\mathcal{C} \cap I_{-n}| = |\mathcal{C} \cap I_{2n}| = 1$ .  $\square$

From this lemma, we deduce that  $|\mathcal{C}| = 2p$  is even. We denote by  $z_{2p}$  the critical point in  $I_{-n}$ ,  $z_1$  the critical point in  $I_{2n}$ , and  $z_2 \geq \dots \geq z_{2p-1}$  the  $2p - 2$  critical points in  $I_{-n+1} \cup \dots \cup I_{2n-1}$ . For  $k = 1, \dots, 2p$ , we define the critical values  $h_k := f(z_k)$ . The next lemma is crucial in establishing the basic properties of  $\rho_{1c}$  (see e.g. Fig. 6).

**Lemma A.3.** (*Orderings of the critical values*) *The critical values are ordered as  $h_1 \geq h_2 \geq \dots \geq h_{2p}$ . Furthermore, there is an absolute constant  $C_0 > 0$  independent of  $\tau$  such that  $h_k \in [-C_0(\tau^{-1}|w|^{-1/2} + |z|) - \sqrt{w}, C_0(\tau^{-1}|w|^{-1/2} + |z|) - \sqrt{w}]$  for  $k = 1, \dots, 2p$ .*

*Proof.* Notice for the equation (2.14), if we multiply both sides with the product of all denominators in  $f$ , we get a polynomial equation  $P_w(m) = 0$  with  $P_w$  being a polynomial of degree  $3n + 1$ . An immediate consequence is that for any fixed  $w > 0$  and  $E \in \mathbb{R}$ ,  $f(\sqrt{w}, m) = E$  can have at most  $3n + 1$  roots in  $m$ . This fact is useful in the proof of this lemma and Lemma 2.3.

For  $i = -n, \dots, 2n$ , define the subset  $J_i(w) := \{m \in I_i : \partial_m f(\sqrt{w}, m) > 0\}$ . From Lemma A.2, we deduce that if  $i = -n + 1, \dots, 2n - 1$ , then  $J_i \neq \emptyset$  if and only if  $I_i$  contains two distinct critical

points of  $f$ , in which case  $J_i$  is an interval. Moreover, we have  $J_{-n} = (-\infty, z_{2p})$  and  $J_{2n} = (z_1, +\infty)$ . Next, we observe that for any  $-n \leq i < j \leq 2n$ , we have  $f(J_i) \cap f(J_j) = \emptyset$ . Otherwise if there were  $E \in f(J_i) \cap f(J_j)$ , we would have  $|\{x : f(x) = E\}| > 3n + 1$ . We hence conclude that the sets  $f(J_i)$ ,  $-n \leq i \leq 2n$  can be strictly ordered. The claim  $h_1 \geq h_2 \geq \dots \geq h_{2p}$  is now reformulated as

$$f(J_i) < f(J_j) \text{ whenever } i < j \text{ and } J_i, J_j \neq \emptyset. \quad (\text{A.17})$$

To prove (A.17), we use a continuity argument. Let  $t \in (0, 1]$  and introduce

$$f^t(m) = -\sqrt{w} + m + w^{-1/2} + \frac{t}{N} \sum_{k=1}^{2n} \frac{C_k^+}{m - x_k} + \frac{t}{N} \sum_{l=1}^n \frac{C_l^-}{m + y_l}.$$

It is easy to check (A.17) holds for small enough  $t > 0$ . We claim that

$$J_i \neq \emptyset \Rightarrow J_i^t \neq \emptyset \text{ for all } t \in (0, 1]. \quad (\text{A.18})$$

This is trivial for  $i = -n, 2n$ . Recall that for  $-n + 1 \leq i \leq 2n - 1$ ,  $J_i^t \neq \emptyset$  is equivalent to  $I_i$  containing two distinct critical points. Moreover,  $\partial_t \partial_m f^t(m) < 0$  in  $I_{-n+1} \cup \dots \cup I_{2n-1}$ , from which we deduce that the number of distinct critical points in each  $I_i$ ,  $i = -n + 1, \dots, 2n - 1$ , does not decrease as  $t$  decreases. This proves (A.18).

Next, suppose that there exist  $i < j$  such that  $J_i, J_j \neq \emptyset$  and  $f(J_i) > f(J_j)$ . From (A.18), we deduce that  $J_i^t, J_j^t \neq \emptyset$  for all  $t \in (0, 1]$ . By a simple continuity argument, we get that  $f^t(J_i^t) > f^t(J_j^t)$  for all  $t \in (0, 1]$ . However, this is impossible for small enough  $t$  as explained before (A.18). This concludes the proof of (A.17).

To prove the second statement of Lemma A.3, we only need to show that  $h_1 \leq C_0(\tau^{-1}|w|^{-1/2} + |z|) - \sqrt{w}$  and  $h_{2p} \geq -C_0(\tau^{-1}|w|^{-1/2} + |z|) - \sqrt{w}$  for some absolute constant  $C_0$ . We only give the proof for  $h_1$ ; the proof for  $h_{2p}$  is similar. At  $z_1$ , we have

$$f(z_1) + \sqrt{w} \leq (z_1 + y_n) \left[ 1 + \frac{1}{N} \sum_{k=1}^{2n} \frac{C_k^+}{(z_1 - x_k)^2} + \frac{1}{N} \sum_{l=1}^n \frac{C_l^-}{(z_1 + y_l)^2} \right] + w^{-1/2} = 2(z_1 + y_n) + w^{-1/2},$$

where we use

$$0 = f'(z_1) = 1 - \frac{1}{N} \sum_{k=1}^{2n} \frac{C_k^+}{(z_1 - x_k)^2} - \frac{1}{N} \sum_{l=1}^n \frac{C_l^-}{(z_1 + y_l)^2}. \quad (\text{A.19})$$

Now we would like to estimate  $z_1 + y_n$ . Again using (A.19), we have that

$$\frac{1}{N} \sum_{k=1}^{2n} \frac{C_k^+}{(z_1 - x_{2n})^2} + \frac{1}{N} \sum_{l=1}^n \frac{C_l^-}{(z_1 - x_{2n})^2} \geq 1.$$

Then by (A.5) we get

$$z_1 - x_{2n} \leq \sqrt{\frac{1}{N} \sum_{k=1}^{2n} C_k^+ + \frac{1}{N} \sum_{l=1}^n C_l^-} \leq \sqrt{5 \frac{\tau^{-1} + |z|^2 + \sqrt{w}|z|}{w}}.$$

Using the above estimates and (A.2)-(A.4), we obtain that

$$f(z_1) \leq 2 \left( \sqrt{5 \frac{\tau^{-1} + |z|^2 + \sqrt{w}|z|}{w}} + \frac{s_1 + |z|^2}{\sqrt{w}} + 2|z| \right) + w^{-1/2} - \sqrt{w} \leq C_0(\tau^{-1}|w|^{-1/2} + |z|) - \sqrt{w}.$$

for some constant  $C_0 > 0$  that does not depend on  $\tau$ .  $\square$

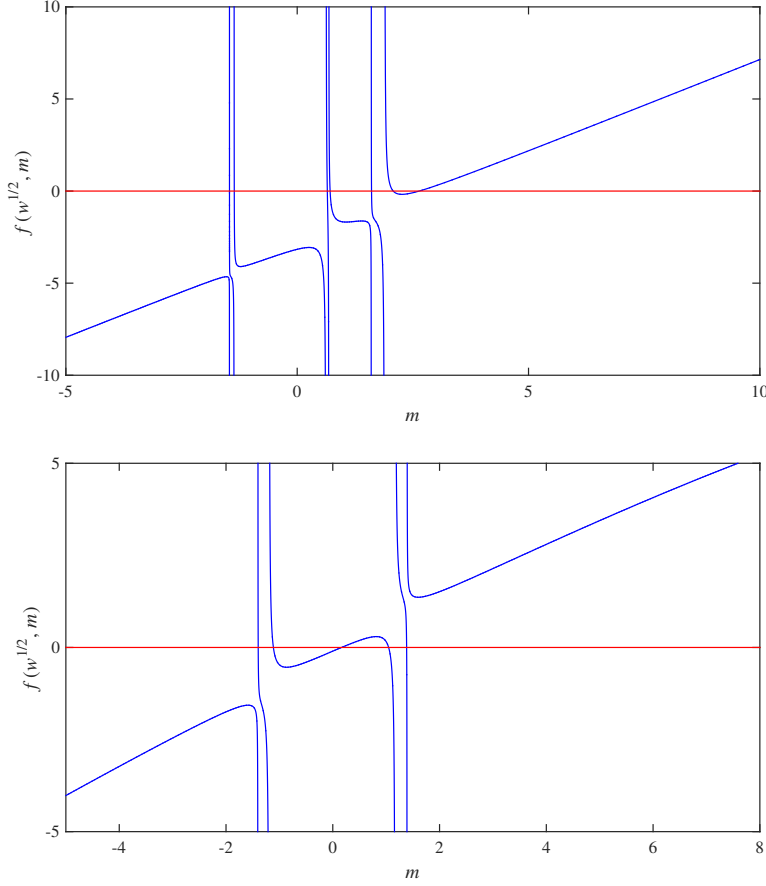


Figure 6: The graphs of  $f(\sqrt{w}, m)$  for the example from Figure 1, i.e.  $\rho_\Sigma = 0.5\delta_{\sqrt{2/17}} + 0.5\delta_{4\sqrt{2/17}}$ . We take  $|z| = 1.5$ , and  $w = 10$  and  $0.01$  in the upper and lower graphs, respectively. In the lower graph, we only plot the five branches near  $m = 0$ . The remaining two branches are far away.

*Proof of Lemma 2.3.* Let  $J(w) := \bigcup_{i=-n}^{2n} J_i(w)$ . Given  $w > 0$  such that  $0 \in f(J(w))$ , then the set  $\{m \in \mathbb{R} : f(\sqrt{w}, m) = 0\}$  has  $3n + 1$  points. Since  $f(\sqrt{w}, m) = 0$  has at most  $3n + 1$  solutions in  $m$ , we deduce that  $m_c(w)$  is real and hence  $m_{1c}(w)$  is also real. Since  $m_{1c}$  is the Stieltjes transform of  $\rho_{1c}$ , we conclude that  $w \notin \text{supp } \rho_{1c}$ . On the other hand, suppose  $w > 0$  and  $0 \notin f(J(w))$ . Then the set of preimages  $\{m \in \mathbb{R} : f(\sqrt{w}, m) = 0\} = \{m \in \mathbb{R} : P_w(m) = 0\}$  has  $3n - 1$  points. Since  $P_w(m)$  is a degree  $3n + 1$  polynomial with real coefficients, we conclude that  $P_w$  has a unique root with positive imaginary part. By the uniqueness of the solution of  $P_{w+i\eta}$  in  $\mathbb{C}_+$  (Lemma 2.2) and the continuity of the roots of  $P_{w+i\eta}$  in  $\eta$ , we conclude that  $\text{Im } m_c(w) > 0$  and  $\text{Im } m_{1c}(w) > 0$  by taking  $\eta \searrow 0$ , i.e.  $w \in \text{supp } \rho_{1c}$ . In sum, we get

$$\text{supp } \rho_{1c} = \overline{\{w > 0 : 0 \notin f(J(w))\}}. \quad (\text{A.20})$$

From Lemma A.3, we see that there exists an absolute constant  $C_1 > 0$  such that if  $w \geq C_1 \tau^{-1}$ , then  $h_1(\omega) \leq C_0(\tau^{-1}|w|^{-1/2} + |z|) - \sqrt{w} < 0$ . Hence fix  $w \geq C_1 \tau^{-1}$ , we have  $0 \in f(J_{2n}(w))$  and

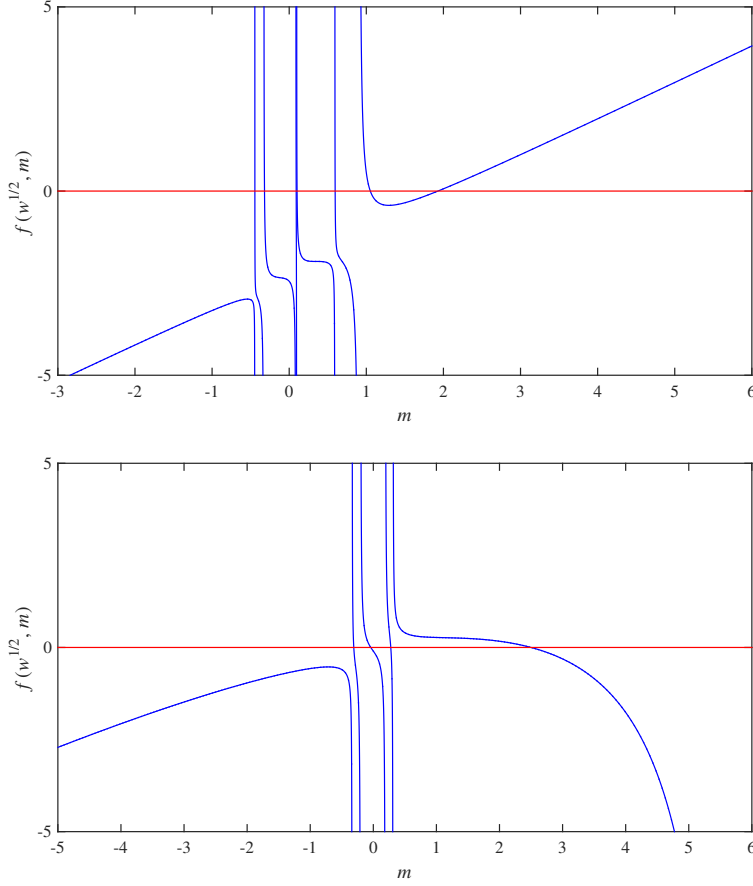


Figure 7: The graphs of  $f(\sqrt{w}, m)$  for the example from Figure 1, i.e.  $\rho_\Sigma = 0.5\delta_{\sqrt{2/17}} + 0.5\delta_{4\sqrt{2/17}}$ . We take  $|z| = 0.5$ , and  $w = 6$  and  $0.01$  in the upper and lower graphs, respectively. In the lower graph, we only plot the five branches near  $m = 0$ . The remaining two branches are far away.

$w \notin \text{supp } \rho_{1c}$  (see the upper graphs in Fig. 6 and 7). This shows that  $\rho_{1c}$  is compactly supported in  $[0, C_1\tau^{-1}]$ . Now we decrease  $w$  so that  $w < s_1 + |z|^2 + 1$ , then using (A.2),

$$h_1(w) > z_1 + w^{-1/2} - \sqrt{w} > \frac{s_1 + |z|^2 + 1 - w}{\sqrt{w}} > 0.$$

By continuity, there must be some  $0 < w < C\tau^{-1}$  such that  $0 \notin f(J(w))$ . Thus  $\text{supp } \rho_{1c} \neq \emptyset$ . By (A.20), it is not hard to see that  $\text{supp } \rho_{1c}$  is a disjoint union of (countably many) closed intervals,

$$\text{supp } \rho_{1c} = \bigcup_k [e_{2k}, e_{2k-1}], \quad (\text{A.21})$$

where  $C_1\tau^{-1} \geq e_1 \geq e_2 \geq \dots$ . Furthermore, for  $e_i$  to be a boundary point, we must have that 0 is a critical value of  $f(\sqrt{e_i}, m)$ , i.e. there is a unique critical point  $m = m_c(e_i)$  such that

$$f(\sqrt{e_i}, m_c(e_i)) = 0, \quad \partial_m f(\sqrt{e_i}, m_c(e_i)) = 0. \quad (\text{A.22})$$

Notice the two equations in (A.22) are equivalent to two polynomial equations in  $(\sqrt{w}, m)$  with order  $3n + 1$  and  $6n$ , respectively. By Bézout's theorem, there are at most finitely many solutions to (A.22). Hence there are finitely many  $e_i$ 's, call them  $e_1 \geq e_2 \geq \dots \geq e_{2L}$ , where  $L \equiv L(n) \in \mathbb{N}$ . To prove the statement about  $e_{2L}$ , we use Lemma A.4 below. This concludes Lemma 2.3.  $\square$

**Lemma A.4.** *If  $1 + \tau \leq |z|^2 \leq 1 + \tau^{-1}$ , there is a constant  $\epsilon(\tau) > 0$  so that  $e_{2L} \geq \epsilon(\tau)$ . If  $|z|^2 \leq 1 - \tau$ ,  $e_{2L} = 0$  and  $\rho_{1c}(x) \sim x^{-1/2}$  when  $x \searrow 0$ .*

*Proof.* By this lemma, the behavior of the leftmost edge  $e_{2L}$  changes essentially when  $z$  crosses the unit circle. From the following proof, we see that the singularity happens at  $|z|^2 = N^{-1} \sum_{i=1}^n l_i s_i$ . Thus the fact that the singular circle has radius 1 comes from our normalization (2.5) for  $T$ .

We first study equation (2.14) when  $w \searrow 0$  in the case  $1 + \tau \leq |z|^2 \leq 1 + \tau^{-1}$ . We calculate the derivative of  $f$  as

$$\begin{aligned} \partial_m f(\sqrt{w}, m) &= 1 + \frac{1}{N} \sum_{i=1}^n l_i s_i \frac{m^2 - |z|^2}{\sqrt{w} m^3 - (s_i + |z|^2) m^2 - \sqrt{w} |z|^2 m + |z|^4} \\ &\quad - \frac{m}{N} \sum_{i=1}^n l_i s_i \frac{\sqrt{w} (m^2 - |z|^2)^2 + 2s_i |z|^2 m}{[\sqrt{w} m^3 - (s_i + |z|^2) m^2 - \sqrt{w} |z|^2 m + |z|^4]^2}. \end{aligned} \quad (\text{A.23})$$

It is easy to see that  $J_0 \neq \emptyset$  for all  $w > 0$ , since  $\partial_m f(\sqrt{w}, 0) = 1 - |z|^{-2} > 0$  (see the lower graph in Fig. 6). Call the end points of  $J_0$  as  $z_k(w) > 0$  and  $z_{k+1}(w) < 0$ . By the definition of  $I_0$ , we have  $z_k < b_1 < |z|$ . Suppose  $z_k = o(|z|)$  as  $w \rightarrow 0$ , then (A.23) gives that  $0 = 1 - |z|^{-2} + o(1)$ , which gives a contradiction. Thus  $z_k \sim |z|$  as  $w \rightarrow 0$ . Now using  $\partial_m f(\sqrt{w}, z_k) = 0$ , we can estimate that

$$\begin{aligned} f(\sqrt{w}, z_k) &= -\sqrt{w} + \frac{z_k^2}{N} \sum_{i=1}^n l_i s_i \frac{\sqrt{w} (z_k^2 - |z|^2)^2 + 2s_i |z|^2 z_k}{[\sqrt{w} z_k^3 - (s_i + |z|^2) z_k^2 - \sqrt{w} |z|^2 z_k + |z|^4]^2} \\ &\geq -\sqrt{w} + \frac{1}{N} \sum_{i=1}^n l_i s_i \frac{2s_i |z|^2 z_k^3}{|z|^8} \geq c - \sqrt{w} \end{aligned} \quad (\text{A.24})$$

for some  $C > 0$  independent of  $w$ , where in the second step we use that

$$\sqrt{w} z_k^3 - (s_i + |z|^2) z_k^2 - \sqrt{w} |z|^2 z_k + |z|^4 > 0, \text{ and } \sqrt{w} z_k^3 - (s_i + |z|^2) z_k^2 - \sqrt{w} |z|^2 z_k < 0$$

which come from that  $0 < z_k < b_i$  for all  $1 \leq i \leq n$ . By (A.24), we can find  $\epsilon$  small enough such that  $f(\sqrt{w}, z_k) > 0$  for all  $0 < w \leq \epsilon$ . In this case  $0 \in f(J_0(w))$  and hence  $w \notin \text{supp } \rho_{1c}$ . In fact, it is not hard to see that there is a solution  $m_0 = \sqrt{w} |z|^2 / (|z|^2 - 1) + o(\sqrt{w}) \in I_0$  such that  $f(\sqrt{w}, m_0) = 0$  and  $\partial_m f(\sqrt{w}, m_0) > 0$ . This proves the first statement of Lemma A.4.

Now we study equation (2.14) when  $|z|^2 \leq 1 - \tau$  and  $w \rightarrow 0$ . For later purpose, we allow  $w$  to be complex and prove a more general result than what we need for this lemma. Let  $w = 0$  in the equation (2.14), we get  $m = 0$  or

$$0 = 1 + \frac{1}{N} \sum_{i=1}^n l_i s_i \frac{m^2 - |z|^2}{-(s_i + |z|^2) m^2 + |z|^4}. \quad (\text{A.25})$$

We define

$$g(x) = 1 + \frac{1}{N} \sum_{i=1}^n l_i s_i \frac{x - |z|^2}{-(s_i + |z|^2)x + |z|^4} = \frac{|z|^2}{N} \sum_{i=1}^n l_i \frac{-x + |z|^2 - s_i}{-(s_i + |z|^2)x + |z|^4}. \quad (\text{A.26})$$

It is easy to see that  $g$  is smooth and decreasing on the intervals defined through

$$K_1 := \left(-\infty, \frac{|z|^4}{s_1 + |z|^2}\right), \quad K_i := \left(\frac{|z|^4}{s_{i-1} + |z|^2}, \frac{|z|^4}{s_i + |z|^2}\right) \quad (i = 2, \dots, n), \quad K_{n+1} := \left(\frac{|z|^4}{s_n + |z|^2}, \infty\right).$$

By the boundary values of  $g$  on these intervals, we see that  $g(x)$  has exactly one zero on intervals  $K_i$  for  $i = 1, \dots, n$ , and has no zero on  $K_{n+1}$ . Since  $g(x) = 0$  is equivalent to a polynomial equation of order  $n$ , it has at most  $n$  solutions. We conclude that all of its solutions are real. Obviously the zeros on the intervals  $K_i$  are positive for  $i = 2, \dots, n$ . Now we study the zero on  $K_1$ . Observe that  $g(0) = 1 - |z|^{-2} < 0$  (as  $|z|^2 \leq 1 - \tau$ ), the zero on  $K_1$  is negative, call it  $-t$ . Moreover, we can verify that  $g(-\tau^{-1}) > 0$  by (A.26), so  $t < \tau^{-1}$ . If  $|z|^2 \geq \tau/2$ , then by the concavity of  $g$  on the  $K_1$ , we get

$$t \geq \frac{g(0)}{g'(0)} \geq \frac{|z|^4(1 - |z|^2)}{s_1} \geq \frac{\tau^4}{4}. \quad (\text{A.27})$$

In the case  $|z|^2 \leq \tau/2$ , we have  $|z|^2 - s_n \leq -\tau/2$  and  $g(|z|^2 - s_n) \leq 0$  by (A.26). Hence we have

$$-t \leq |z|^2 - s_n \leq -\tau/2. \quad (\text{A.28})$$

Combining (A.27) and (A.28), we get that  $c\tau^4 \leq t \leq \tau^{-1}$  for some constant  $c > 0$ .

Now we return to the self-consistent equation (2.14). The previous discussions show that

$$f(0, i\sqrt{t}) = 0, \quad t \geq c\tau^4.$$

It is easy to see that there exists constants  $c_1, \tau' > 0$  such that

$$|-(s_i + |z|^2)m^2 + |z|^4 + \sqrt{w}(m^3 - |z|^2m)| \geq c_1 \text{ for } |m - i\sqrt{t}| \leq \tau'. \quad (\text{A.29})$$

First we consider the case  $|z| \geq \epsilon > 0$ . Expanding  $f(\sqrt{w}, m)$  around  $(0, i\sqrt{t})$  and using (A.29),

$$0 = \partial_{\sqrt{w}}f(0, i\sqrt{t})\sqrt{w} + \partial_m f(0, i\sqrt{t})(m - i\sqrt{t}) + o(\sqrt{w}) + o(m - i\sqrt{t}). \quad (\text{A.30})$$

By (A.23),

$$\partial_{\sqrt{w}}f(\sqrt{w}, m) = -1 - \frac{m^2}{N} \sum_{i=1}^n l_i s_i \frac{(m^2 - |z|^2)^2}{[-(s_i + |z|^2)m^2 + |z|^4 + \sqrt{w}(m^3 - |z|^2m)]^2}, \quad (\text{A.31})$$

and (A.29), we get  $|\partial_{\sqrt{w}}f(0, i\sqrt{t})| \leq C$  and

$$\partial_m f(0, i\sqrt{t}) = \frac{t}{N} \sum_{i=1}^n l_i s_i \frac{2s_i |z|^2}{[(s_i + |z|^2)t + |z|^4]^2} \geq c_2 \quad (\text{A.32})$$

for some  $c_2 > 0$ . Using (A.32), we get from (A.30) that

$$m - i\sqrt{t} = O(\sqrt{w}), \quad \text{if } |z| \geq \epsilon. \quad (\text{A.33})$$

In particular, this shows that  $|m| \approx \text{Im } m \sim 1$  as  $w \rightarrow 0$ .

Then assume that  $|z|^2 < \epsilon$ , for sufficiently small  $\epsilon$ . From  $g(-t) = 0$  and (A.26), we get that

$$\frac{1}{N} \sum_{i=1}^n l_i \frac{t + |z|^2 - s_i}{(s_i + |z|^2)t + |z|^4} = 0. \quad (\text{A.34})$$



From the leading order term, we get  $t^{-1} = t_0^{-1} + O(|z|^2)$ , where  $t_0 := (N^{-1} \sum_i l_i/s_i)^{-1}$ . Expanding (A.34) to the first order term of  $|z|^2$ , we get

$$t = t_0 + \left( \frac{t_0^2}{N} \sum_i \frac{l_i}{s_i^2} - 2 \right) |z|^2 + O(|z|^4). \quad (\text{A.35})$$

Now we write equation (2.14) as

$$F(\sqrt{w}, m) = 0, \quad (\text{A.36})$$

where  $F(\sqrt{w}, m) := f(\sqrt{w}, m)/m$ . Expanding  $F$  around  $(0, i\sqrt{t})$  and using (A.29), we get

$$\begin{aligned} 0 = & \partial_{\sqrt{w}} F(0, i\sqrt{t}) \sqrt{w} + \partial_m F(0, i\sqrt{t}) (m - i\sqrt{t}) + \partial_m \partial_{\sqrt{w}} F(0, i\sqrt{t}) (m - i\sqrt{t}) \sqrt{w} \\ & + \frac{1}{2} \partial_{\sqrt{w}}^2 F(0, i\sqrt{t}) w + \frac{1}{2} \partial_m^2 F(0, i\sqrt{t}) (m - i\sqrt{t})^2 + o(w, |m - i\sqrt{t}|^2, |m - i\sqrt{t}| \sqrt{w}). \end{aligned} \quad (\text{A.37})$$

We can calculate that (the partial derivatives of  $F$  can be obtained using (A.23) and (A.31))

$$\partial_m F(\sqrt{w}, i\sqrt{t}) = -\frac{2i|z|^2 + 2\sqrt{wt_0}}{t_0^{3/2}} + o(|z|^2, \sqrt{w}), \quad (\text{A.38})$$

$$\partial_{\sqrt{w}} F(\sqrt{w}, i\sqrt{t}) = (i|z|^2 + 2\sqrt{wt_0}) \frac{\sqrt{t_0}}{N} \sum_{j=1}^n \frac{l_j}{s_j^2} + o(|z|^2, \sqrt{w}). \quad (\text{A.39})$$

From (A.38) and (A.39), we get that

$$\begin{aligned} \partial_m F(0, i\sqrt{t}) &= -\frac{2i|z|^2}{t_0^{3/2}} + o(|z|^2), \quad \partial_{\sqrt{w}} F(0, i\sqrt{t}) = \frac{i|z|^2 \sqrt{t_0}}{N} \sum_{j=1}^n \frac{l_j}{s_j^2} + o(|z|^2), \\ \partial_m \partial_{\sqrt{w}} F(0, i\sqrt{t}) &= -\frac{2}{t_0} + O(|z|^2), \quad \partial_{\sqrt{w}}^2 F(0, i\sqrt{t}) = \frac{2t_0}{N} \sum_{j=1}^n \frac{l_j}{s_j^2} + O(|z|^2), \quad \partial_m^2 F(0, i\sqrt{t}) = O(|z|^2). \end{aligned}$$

Plugging the above results into (A.37), we get that

$$\begin{aligned} 0 = & \left[ \frac{i|z|^2 \sqrt{t_0} + \sqrt{wt_0}}{N} \sum_{j=1}^n \frac{l_j}{s_j^2} + o(|z|^2) \right] \sqrt{w} + \left[ -2 \frac{i|z|^2 + \sqrt{wt_0}}{t_0^{3/2}} + o(|z|^2) \right] (m - i\sqrt{t}) \\ & + o(w, |m - i\sqrt{t}|^2, |m - i\sqrt{t}| \sqrt{w}). \end{aligned} \quad (\text{A.40})$$

Observing that  $|i|z|^2 \sqrt{t_0} + \sqrt{wt_0}| \sim |z|^2 + \sqrt{|w|}$ , we get

$$m - i\sqrt{t} = \left[ \frac{t_0^2}{2N} \sum_{j=1}^n \frac{l_j}{s_j^2} + O(|w|^{1/2} + |z|^2) \right] \sqrt{w}, \quad \text{if } |z| < \epsilon. \quad (\text{A.41})$$

Combing (A.33) and (A.41), we get that if  $|z|^2 < 1 - \tau$ ,  $m = i\sqrt{t} + O(\sqrt{w})$  when  $w \rightarrow 0$ . In particular, this shows that  $|m| \approx \text{Im } m \sim 1$  when  $w \rightarrow 0$ . Finally we conclude the proof of Lemma A.4 by using that  $m_{1c}(w) = m_c(w)w^{-1/2} - 1$ .  $\square$

To prove Proposition 2.14, we need the following lemma, which is a consequence of the edge regularity conditions (2.18) and (2.19).

**Lemma A.5.** Suppose  $e_k \neq 0$  is a regular edge. Then  $|m_{1c}(w) - m_{1c}(e_k)| \sim |w - e_k|^{1/2}$  as  $w \rightarrow e_k$  and  $\min_{l \neq k} |e_l - e_k| \geq \delta$  for some constant  $\delta > 0$ .

*Proof.* Denote  $m_k := m_c(e_k)$  and let  $w \rightarrow e_k$ . Notice by Lemma 2.3, if  $e_k \neq 0$ , we have

$$\epsilon \leq e_k \leq C\tau^{-1}. \quad (\text{A.42})$$

Then we expand  $f$  around  $(\sqrt{e_k}, m_k)$  to get that

$$\begin{aligned} 0 = & \partial_{\sqrt{w}} f(\sqrt{e_k}, m_k)(\sqrt{w} - \sqrt{e_k}) + \frac{1}{2} \partial_m^2 f(\sqrt{e_k}, m_k)(m_c(w) - m_k)^2 \\ & + O[|\sqrt{w} - \sqrt{e_k}|^2 + |m_c(w) - m_k|^3 + |\sqrt{w} - \sqrt{e_k}||m_c(w) - m_k|], \end{aligned} \quad (\text{A.43})$$

where by (A.31),

$$\partial_{\sqrt{w}} f(\sqrt{e_k}, m_k) = -1 - \frac{m_k^2}{N} \sum_{i=1}^n l_i s_i \frac{(m_k^2 - |z|^2)^2}{e_k(m_k - a_i)^2(m_k - b_i)^2(m_k + c_i)^2}, \quad (\text{A.44})$$

and by (A.1),

$$\partial_m^2 f(\sqrt{e_k}, m_k) = \frac{2}{N} \sum_{i=1}^n l_i s_i \left[ \frac{A_i}{(m_k - a_i)^3} + \frac{B_i}{(m_k - b_i)^3} + \frac{C_i}{(m_k + c_i)^3} \right], \quad (\text{A.45})$$

Applying (A.2)-(A.5), (A.42) and the conditions (2.18)-(2.19) to (A.44) and (A.45), we get that

$$1 \leq |\partial_{\sqrt{w}} f(\sqrt{e_k}, m_k)| \leq C_1, \quad \epsilon \leq |\partial_m^2 f(\sqrt{e_k}, m_k)| \leq C_2 \quad (\text{A.46})$$

for some  $C_1, C_2 > 0$ . Similarly, if  $|w - e_k| \leq \tau'$  and  $|m_c(w) - m_k| \leq \tau'$  for some sufficiently small  $\tau'$ , using the condition (2.18) we can get that

$$\max \left\{ |\partial_m^3 f(\sqrt{w}, m_c(w))|, |\partial_{\sqrt{w}}^2 f(\sqrt{w}, m_c(w))|, |\partial_m \partial_{\sqrt{w}} f(\sqrt{w}, m_c(w))| \right\} \leq C_3. \quad (\text{A.47})$$

Plug them into equation (A.43), for  $|w - e_k| \leq \tau'$  and  $|m_c(w) - m_k| \leq \tau'$ , we get  $|m_c(w) - m_k| \sim |\sqrt{w} - \sqrt{e_k}|^{1/2}$  and

$$-\partial_{\sqrt{w}} f(\sqrt{e_k}, m_k)(\sqrt{w} - \sqrt{e_k}) + O(|\sqrt{w} - \sqrt{e_k}|^{3/2}) = \frac{1}{2} \partial_m^2 f(\sqrt{e_k}, m_k)(m_c(w) - m_k)^2. \quad (\text{A.48})$$

By (A.42), we immediately get that  $|\sqrt{w} - \sqrt{e_k}| \sim |w - e_k|$  and  $|m_c(w) - m_k| \sim |m_{1c}(w) - m_{1c}(e_k)|$ , which proves the first part of the lemma. By (A.48), if  $w$  is real and  $|w - e_k| \leq \tau'$ , we have that

$$m_c(w) - m_k = \left[ \frac{-2\partial_{\sqrt{w}} f(\sqrt{e_k}, m_k)}{\partial_m^2 f(\sqrt{e_k}, m_k)} + O(|\sqrt{w} - \sqrt{e_k}|^{1/2}) \right]^{1/2} (\sqrt{w} - \sqrt{e_k})^{1/2}. \quad (\text{A.49})$$

Thus on a sufficiently small interval  $U = [e_k - \delta, e_k + \delta]$ ,  $m_c(w)$  has positive imaginary part for  $w$  on one side of  $e_k$  and  $m_c(w)$  is real for  $w$  on the other side. Hence  $U$  does not contain another edge. This shows that  $\min_{l \neq k} |e_l - e_k| \geq \delta$ .  $\square$

*Proof of Proposition 2.14.* The properties of  $\rho_{1c}$  have been proved in Lemmas 2.3, A.4 and A.5, and included in the Definition 2.4. Since  $\text{supp } \rho_{2c} = \text{supp } \rho_{1c}$  by the discussions after Lemma 2.2, we

immediately get property (i) for  $\rho_{2c}$ . The conclusion  $\rho_{2c}$  being a probability measure is due to the definition of  $m_2$  in (2.34) and the fact that  $m_{2c}$  is the almost sure limit of  $m_2$ .

The properties (ii) and (iv) for  $\rho_{2c}$  can be easily obtained by plugging  $m_{1c}$  into (2.9). To prove the property (iii) for  $\rho_{2c}$ , we need to know the behavior of  $\text{Im } m_{2c}(w)$  when  $w \rightarrow e_j$  along the real line. By (2.9), it suffices to prove that if  $|x - e_j| \leq \tau'$  for some small enough  $\tau' > 0$ , then

$$|-w(1 + m_{1c})^2 + |z|^2| = |m_c^2 - |z|^2| \geq \epsilon$$

for some constant  $\epsilon > 0$ . Suppose that  $|m_c^2(w) - |z|^2| = o(1)$ . Plugging  $m_c$  into  $\partial_m f(\sqrt{w}, m_c)$  in (A.23), and using condition (2.18) and Lemma A.5, we get that

$$\partial_m f(\sqrt{w}, m_c(w)) = -1 + O(|m_c^2 - |z|^2|). \quad (\text{A.50})$$

Again using condition (2.18) and Lemma A.5, we can bound  $\partial_{\sqrt{w}} \partial_m f(\sqrt{w}, m_c(w))$  and  $\partial_m^2 f(\sqrt{w}, m_c(w))$  for  $w$  near  $e_j$ . Thus we shall have that

$$0 = \partial_m f(\sqrt{e_j}, m_c(e_j)) = \partial_m f(\sqrt{w}, m_c(w)) + O(|w - e_j|^{1/2}) = -1 + O(|m_c^2 - |z|^2| + |w - e_j|^{1/2}). \quad (\text{A.51})$$

This gives a contradiction. Thus we must have a lower bound for  $|m_c^2 - |z|^2|$ .  $\square$

*Remark:* Here we add a small remark on Example 2.8. Given the assumptions in Example 2.8, it is easy to see that  $f$  can only take critical values on intervals  $I_{-n}$ ,  $I_0$ ,  $I_n$  and  $I_{2n}$ , since  $\max\{|a_i - a_{i-1}|, |b_i - b_{i-1}|, |c_i - c_{i-1}|\} \rightarrow 0$  in this case. Thus the number of connected components of  $\text{supp } \rho_{1c}$  is independent of  $n$ , and all the edges and the bulk components are regular as in Example 2.7.

## A.2 Proof of Lemmas 3.7 and 3.8

We first prove Lemma 3.7. We consider the five cases separately.

*Case 1:* For  $w = E + i\eta \in \mathbf{D}_k^b(\zeta, \tau', N)$ , we have

$$m_{1c}(w) = \int_{\mathbb{R}} \frac{\rho_{1c}(x)}{x - (E + i\eta)} dx, \quad \text{Im } m_{1c}(w) = \int_{\mathbb{R}} \frac{\rho_{1c}(x, z)\eta}{(x - E)^2 + \eta^2} dx. \quad (\text{A.52})$$

By the regularity condition of Definition 2.4 (ii), we get immediately  $\text{Im } m_{1c} \sim 1$ . Since  $\text{Im } m_{1c} \leq |1 + m_{1c}| \leq C$  by Proposition 2.15, we get  $|1 + m_{1c}| \sim 1$ . Notice  $wm_{1c}$  can be expressed as

$$wm_{1c}(w) = \int_{\mathbb{R}} \frac{w\rho_{1c}(x, z)}{x - w} dx = - \int_{\mathbb{R}} \rho_{1c}(x, z) dx + \int_{\mathbb{R}} \frac{x\rho_{1c}(x, z)}{x - w} dx.$$

By the same argument as above and using the fact that  $x \geq \tau'$  for  $x \in [e_{2k} + \tau', e_{2k-1} - \tau']$ , we get

$$\text{Im}(wm_{1c}) = \text{Im} \int_{\mathbb{R}} \frac{x\rho_{1c}(x, z)}{x - w} dx \sim 1.$$

Since the imaginary parts of  $-w$  and  $|z|^2/(1 + m_{1c})$  are both negative, we get

$$\text{Im} \left[ -w(1 + m_{1c}) + \frac{|z|^2}{1 + m_{1c}} \right] \leq -\text{Im}(wm_{1c}). \quad (\text{A.53})$$

Using the bounds for  $m_{1c}$  and  $\text{Im } m_{1c}$  proved above, it is easy to see

$$\left| -w(1 + m_{1c}) + \frac{|z|^2}{1 + m_{1c}} \right| = O(1). \quad (\text{A.54})$$

Equations (A.53) and (A.54) together give that  $\text{Im } m_{2c} \sim 1$  and  $|m_{2c}| \sim 1$ . Similarly, we can also prove that

$$wm_{2c} = \left[ -(1 + m_{1c}) + \frac{|z|^2}{w(1 + m_{1c})} \right]^{-1} \in \mathbb{C}_+$$

and  $\text{Im}(wm_{2c}) \sim 1$ . Now (3.29) follows from

$$\text{Im} \left( w + s_i wm_{2c} - \frac{|z|^2}{1 + m_{1c}} \right) \geq s_i \text{Im}(wm_{2c}).$$

*Case 2:* For  $w = E + i\eta \in \mathbf{D}^o(\zeta, \tau', N)$ , using (A.52) and  $\text{dist}(E, \text{supp } \rho_{1,2c}) \geq \tau'$ , we immediately get  $\text{Im } m_{1,2c} \sim \eta$ . Now we prove the other estimates.

We first prove (3.29). If  $\eta \sim 1$ , the proof is exactly the same as in Case 1. Hence we assume  $\eta \leq c'$ , where  $c' \equiv c'(\tau, \tau') > 0$  is sufficiently small. We separate it into two cases.

(i) Suppose  $E \sim 1$ . We shall prove that

$$\min_i \{|m_c(w) - a_i(w)|, |m_c(w) - b_i(w)|, |m_c(w) + c_i(w)|\} \geq \epsilon', \quad (\text{A.55})$$

for some constant  $\epsilon'$ . This leads immediately to (3.29) since

$$\left| w \left( 1 + s_i \frac{1 + m_{1c}}{-w(1 + m_{1c})^2 + |z|^2} \right) (1 + m_{1c}) - |z|^2 \right| = \left| \frac{\sqrt{w}(m_c - a_i)(m_c - b_i)(m_c + c_i)}{-m_c^2 + |z|^2} \right|. \quad (\text{A.56})$$

For  $p_i = \sqrt{E}m^3 - (s_i + |z|^2)m^2 - \sqrt{E}|z|^2m + |z|^4$ , it is not hard to prove that its roots  $a_i(E)$ ,  $b_i(E)$  and  $-c_i(E)$  decrease as  $E$  increase. Since  $E \notin \text{supp } \rho_{1c}$ , we have  $m_{1c}(E) \in \mathbb{R}$  and

$$\frac{dm_{1c}(E)}{dE} = \int_{\mathbb{R}} \frac{\rho_{1c}(x, z)}{(x - E)^2} dx \geq 0.$$

So  $m_{1c}(E)$  (and hence  $m_c(E)$ ) increases as  $E$  increases. If  $e_k$  is the smallest edge that is bigger than  $E$ , then for  $a_i(E)$  bigger than  $m_c(E)$ , we have that

$$a_i(E) - m_c(E) \geq a_i(e_k) - m_c(e_k) + \epsilon(\tau') \geq \epsilon(\tau'), \quad (\text{A.57})$$

by using  $|E - e_k| \geq \tau'$  (see (2.42)). On the other hand, If  $e_{k-1}$  is the largest edge value that is smaller than  $E$ , then for  $a_i(E)$  smaller than  $m_c(E)$ , we have that

$$m_c(E) - a_i(E) \geq m_c(e_{k-1}) - a_i(e_{k-1}) + \epsilon(\tau') \geq \epsilon(\tau'). \quad (\text{A.58})$$

Applying the same arguments to  $b_i(E)$  and  $-c_i(E)$ , we get

$$\min_i \{|m_c(E) - a_i(E)|, |m_c(E) - b_i(E)|, |m_c(E) + c_i(E)|\} \geq \epsilon \quad (\text{A.59})$$

for  $E \in (e_{2k-1}, e_{2k})$  for some  $k$ . Now we are only left with the case  $E < e_{2L}$ , the rightmost edge, when  $|z|^2 \geq 1 + \tau$ . In this case, we have seen that  $0 < m_c(E) < b_i(E)$  for all  $i$  in the proof of Lemma A.4. Thus we can use (A.57) to get lower bounds for  $|m_c(E) - a_i(E)|$  and  $|m_c(E) - b_i(E)|$ . Since  $c_i(E) \sim 1$  in this case (e.g. by (A.4) and using  $E, |z| \sim 1$ ),  $|m_c(E) + c_i(E)| \geq \epsilon$  is trivial. Again we get the estimate (A.59).

Then we consider  $w = E + i\eta$  with  $\eta \leq c'$ . First it is easy to check that  $a_i(E + i\eta)$ ,  $b_i(E + i\eta)$  and  $c_i(E + i\eta)$  are continuous in  $\eta$ . On the other hand for  $m_c(E + i\eta)$ , we have

$$\partial_w m_{1c}(w) = \int_{\mathbb{R}} \frac{\rho_{1c}(x, z)}{(x - w)^2} dx \leq C \quad (\text{A.60})$$

by the condition  $\text{dist}(E, \text{supp } \rho_{1c}) \geq \tau'$ . Thus we immediately get  $|m_c(E + i\eta) - m_c(E)| = O(\eta)$ . Hence as long as  $c'$  is small enough, (A.55) is true, which further gives (3.29).

(ii) Suppose  $w = E + i\eta \rightarrow 0$ , in which case we must have  $|z|^2 \geq 1 + \tau$  and  $E < e_{2L}$ . Using  $|m_{1,2c}(w)| \sim 1$  by Proposition 2.15, we can calculate directly that

$$\left| w(1 + s_i m_{2c})(1 + m_{1c}) - |z|^2 \right| = \left| |z|^2 + O(w) \right| \geq c.$$

This concludes the proof of (3.29).

Then we show that  $|1 + m_{1c}| \sim 1$  for  $w \in \mathbf{D}^o$  and  $\eta \leq c'$ . We again divide it into two cases. First suppose  $|w| \sim 1$ . If  $|m_c|$  can be arbitrarily small, then by (3.29) we get that

$$f(\sqrt{w}, m_c) = -\sqrt{w} + O(m_c) \neq 0,$$

which gives a contradiction. Then suppose  $w = E + i\eta \rightarrow 0$  when  $|z|^2 \geq 1 + \tau$  and  $E < e_{2L}$ . We have seen in the proof of Lemma A.4 that

$$m_c(E) = \sqrt{E} \frac{|z|^2}{|z|^2 - 1} + o\left(\sqrt{E}\right) \Rightarrow 1 + m_{1c}(E) = \frac{|z|^2}{|z|^2 - 1} + o(1).$$

Then using (A.60), we get

$$|1 + m_{1c}(E + i\eta)| = \left| \frac{|z|^2}{|z|^2 - 1} + o(1) + O(\eta) \right| \sim 1.$$

Finally we have  $|m_{2c}| \sim 1$  for  $w \in \mathbf{D}^o$  and  $\eta \leq c'$  by Proposition 2.15.

*Case 3:* For regular edge  $e_k \neq 0$ , we always have  $e_k \geq \epsilon$  for some  $\epsilon > 0$  by Lemma A.4. Thus we always have  $|w| \sim 1$  for  $w = E + i\eta \in \mathbf{D}_k^e(\zeta, \tau', N)$  as long as  $\tau'$  is sufficiently small. If  $\eta \sim 1$ , then  $\sqrt{\kappa} + \eta \sim \eta/\sqrt{\kappa} + \eta \sim 1$  and the proof is exactly the same as in Case 1. Now we pick  $\tau'$  small and consider the case  $\eta \leq \tau'$ . By the regularity assumption (2.18) and Lemma A.5, we have

$$\min_{1 \leq i \leq n} \{|m_c(w) - a_i(w)|, |m_c(w) - b_i(w)|, |m_c(w) + c_i(w)|\} \geq \epsilon/2 \quad (\text{A.61})$$

uniformly in  $w \in \{w \in \mathbf{D}_k^e(\zeta, \tau', N) : \kappa(w) + \eta(w) \leq 2\tau'\}$ , provided  $\tau'$  is sufficiently small. The above bound implies (3.29). If  $m_c(w) \rightarrow 0$ , then using (3.29) we get from  $f(\sqrt{w}, m_c) = 0$  that  $-\sqrt{w} + O(m_c) = 0$ , which gives a contradiction. Thus we must have  $|1 + m_{1c}| \sim |m_c| \sim 1$ . To show  $|m_{2c}| \sim 1$ , we can use Proposition 2.15.

We still need to prove the estimates for  $\text{Im } m_{1,2c}$  when  $\eta \leq \tau'$ . Recall the expansion (A.48) around  $e_k$  and equation (A.49). Notice both  $\partial_{\sqrt{w}} f(\sqrt{e_k}, m_k)$  and  $\partial_m^2 f(\sqrt{e_k}, m_k)$  are real (as  $e_k$  and  $m_k$  are real). Suppose  $k$  is odd, then  $\text{Im } m_c(E) = 0$  for  $E \searrow e_k$  (i.e.  $E \notin \text{supp } \rho_c$ ) and  $\text{Im } m_c(E) > 0$  for  $E \nearrow e_k$  (i.e.  $E \in \text{supp } \rho_c$ ). Thus (A.49) gives

$$m_c(w) - m_k = C_k(w)(w - e_k)^{1/2} + D_k(w),$$

with  $C_k > 0$ ,  $C_k \sim 1$ ,  $|D_k| = O(|w - e_k|)$  and  $\text{Im } D_k \sim \eta$ . Then for  $E \geq e_k$ , we have

$$\text{Im } m_c(E + i\eta) \sim \text{Im}(\kappa + i\eta)^{1/2} + O(\eta) \sim \frac{\eta}{\sqrt{\kappa + \eta}},$$

and for  $E \leq e_k$ , we have

$$\text{Im } m_c(E + i\eta) \sim \text{Im}(-\kappa + i\eta)^{1/2} + O(\eta) \sim \sqrt{\kappa + \eta}.$$

If  $k$  is even, the proof is the same except that in this case

$$m_c(w) - m_k = C_k(w)(e_k - w)^{1/2} + D_k(w).$$

For  $m_{1c}(w)$  and  $m_{2c}(w)$ , we get the conclusion by noticing  $w \approx e_k$  and

$$\text{Im } m_{1c} = \text{Im} \left( w^{-1/2} m_c \right) \sim \text{Im } m_c(w), \quad \text{Im } m_{2c} = \text{Im} \left[ \frac{m_c}{\sqrt{w}(-m_c^2 + |z|^2)} \right] \sim \text{Im } m_c(w).$$

*Case 4:* Again if  $\eta \sim 1$ , the proof is the same as in Case 1. If  $|w| \leq 2\tau'$  for small enough  $\tau'$ , in the proof of Lemma A.4, we have seen that  $m_c = i\sqrt{t} + O(\sqrt{w})$ , which gives the first equation in (3.26). Plugging it into (2.9), we get the second equation in (3.26). Taking the imaginary part, we obtain (3.27). Finally using (3.26), we get (3.29) easily.

*Case 5:* For  $w = E + i\eta \in \mathbf{D}_L(\zeta, N)$ , the bounds for  $m_{1,2}$  and  $\text{Im } m_{1,2}$  in (3.28) follows from (A.52) directly.

Finally we prove Lemma 3.8. The estimates (3.31) and (3.32) follow immediately from (2.32), (3.29) and (3.30). For (3.33), we can write

$$\Pi_{\mathbf{v}\mathbf{v}} = \left\langle \mathbf{v}, \begin{pmatrix} U & 0 \\ 0 & U \end{pmatrix} \Pi_d \begin{pmatrix} U^\dagger & 0 \\ 0 & U^\dagger \end{pmatrix} \mathbf{v} \right\rangle = (\Pi_d)_{\mathbf{u}\mathbf{u}} = \sum_{i=1}^N \langle u_{[i]}, \pi_{[i]c} u_{[i]} \rangle,$$

where

$$\mathbf{u} := \begin{pmatrix} U^\dagger & 0 \\ 0 & U^\dagger \end{pmatrix} \mathbf{v}, \quad u_{[i]} := \begin{pmatrix} u_i \\ u_{\bar{i}} \end{pmatrix}.$$

To control  $\text{Im } \Pi_{\mathbf{v}\mathbf{v}}$ , it is enough to bound  $\langle u_{[i]}, \pi_{[i]c} u_{[i]} \rangle$  for each  $i$ .

We first consider Cases 1-4 of Lemma 3.7. By the definition of  $\pi_{[i]c}$  in (2.32), we get

$$\begin{aligned} \text{Im } \pi_{ii,c} &= |u_i|^2 \text{Im} \left[ -w(1 + |d_i|^2 m_{2c}) + \frac{|z|^2}{1 + m_{1c}} \right]^{-1} \leq \frac{C}{|w|} \text{Im} \left[ w(1 + |d_i|^2 m_{2c}) - \frac{|z|^2}{1 + m_{1c}} \right] \\ &= \frac{C}{|w|} \left[ (1 + |d_i|^2 \text{Re } m_{2c}) \text{Im } w + |d_i|^2 (\text{Re } w) \text{Im } m_{2c} + \frac{|z|^2}{|1 + m_{1c}|^2} \text{Im } m_{1c} \right], \end{aligned}$$

where in the second step we use (3.29) and  $|1 + m_{1c}| \sim |w|^{-1/2}$ . In the first three cases of Lemma 3.7, we have  $|w| \sim 1$  and  $\text{Im } w = O(\text{Im } m_{1c})$ , which give that  $\text{Im } \pi_{ii,c} \leq C \text{Im}(m_{1c} + m_{2c})$ . In case 4 of Lemma 3.7, we use  $|\text{Im } w| + |\text{Re } w| + |1 + m_{1c}|^{-2} = O(|w|)$  and  $\text{Im } m_{1,2c} \sim |w|^{-1/2}$  to get that  $\text{Im } \pi_{ii,c} \leq C \text{Im}(m_{1c} + m_{2c})$ . Similarly we have the bound  $\text{Im } \pi_{\bar{i}\bar{i},c} \leq C \text{Im}(m_{1c} + m_{2c})$ . Finally we can estimate the following term using similar methods,

$$\begin{aligned} \text{Im} (\bar{u}_{\bar{i}} u_i \pi_{\bar{i}i,c} + \bar{u}_i u_{\bar{i}} \pi_{i\bar{i},c}) &= 2 \text{Re} (\bar{u}_i u_{\bar{i}} z) \text{Im} \left\{ w^{-1/2} [w(1 + |d_i|^2 m_{2c})(1 + m_{1c}) - |z|^2]^{-1} \right\} \\ &\leq C \text{Re} (\bar{u}_i u_{\bar{i}} z) \text{Im}(m_{1c} + m_{2c}) \leq C (|u_i|^2 + |u_{\bar{i}}|^2) \text{Im}(m_{1c} + m_{2c}). \end{aligned}$$

Combining the above estimates we get  $\text{Im} \langle u_{[i]}, \pi_{[i]c} u_{[i]} \rangle \leq C |u_{[i]}|^2 \text{Im}(m_{1c} + m_{2c})$ , which implies (3.33). For the Case 5 of Lemma 3.7, we use (3.28) and (3.32) to get

$$\text{Im} \langle u_{[i]}, \pi_{[i]c} u_{[i]} \rangle \leq |u_{[i]}|^2 \|\pi_{[i]c}\| \leq C |u_{[i]}|^2 \text{Im}(m_{1c} + m_{2c}).$$

### A.3 Proof of Lemma 3.10 and Lemma 2.2

We first prove Lemma 3.10. During the proof, we also use the following equivalent definition of the stability expressed in terms of  $m = \sqrt{w}(1 + m_1)$ ,  $u = \sqrt{w}(1 + u_1)$  and  $f(\sqrt{w}, m)$ . Suppose the assumptions in Definition 3.9 holds. Let  $w \in \mathbf{D}$  and suppose that for all  $w' \in L(w)$  we have  $|f(\sqrt{w}, u)| \leq |w|^{1/2}\delta(w)$ . Then

$$|u(w) - m_c(w)| \leq \frac{C|w|^{1/2}\delta}{\sqrt{\kappa} + \eta + \delta}. \quad (\text{A.62})$$

*Case 1:* We take over the notations in Definition 3.9 and abbreviate  $R := f(\sqrt{w}, u)$ , so that  $|R| \leq |w|^{1/2}\delta$ . Then we write the equation  $f(\sqrt{w}, u) - f(\sqrt{w}, m_c) = R$  as

$$\alpha(u - m_c)^2 + \beta(u - m_c) = R, \quad (\text{A.63})$$

where using (A.1),  $\alpha$  and  $\beta$  can be expressed as

$$\alpha := \frac{1}{N} \sum_{i=1}^n l_i s_i \left[ \frac{A_i}{(u - a_i)(m_c - a_i)^2} + \frac{B_i}{(u - b_i)(m_c - b_i)^2} + \frac{C_i}{(u + c_i)(m_c + c_i)^2} \right], \quad (\text{A.64})$$

and

$$\beta := 1 - \frac{1}{N} \sum_{i=1}^n l_i s_i \left[ \frac{A_i}{(m_c - a_i)^2} + \frac{B_i}{(m_c - b_i)^2} + \frac{C_i}{(m_c + c_i)^2} \right] = \partial_m f(\sqrt{w}, m_c). \quad (\text{A.65})$$

We shall prove that

$$|\alpha| + |\partial_u \alpha| \leq C, \quad |\beta| \sim 1, \quad (\text{A.66})$$

for  $w \in \mathbf{D}_k^b$  and  $u$  satisfying  $|u - m_c| \leq (\log N)^{-1/3}$ . If  $|u - m_c| \leq (\log N)^{-1/3}$ , we also have  $\text{Im } u \sim 1$ . By (3.29),

$$\min_i \{|m_c - a_i|, |m_c - b_i|, |m_c + c_i|\} \geq \epsilon \quad (\text{A.67})$$

for some  $\epsilon > 0$ . Replacing the  $m_c$  in (3.29) with  $u$ , we also get that

$$\min_i \{|u - a_i|, |u - b_i|, |u + c_i|\} \geq \epsilon' \quad (\text{A.68})$$

for some  $\epsilon' > 0$ . Using (A.67) and (A.68), we get immediately that  $|\alpha| + |\partial_u \alpha| + |\beta| \leq C$ . What remains is the proof of the lower bound  $|\beta| \geq c$ . If  $\text{Im } w \geq \epsilon$  for some constant  $\epsilon > 0$ , the lower bound follows from Lemma A.6 below. If  $\text{Im } w \leq \epsilon$  for a sufficiently small  $\epsilon$ , the lower bound follows from Lemma A.7 below. Now given the bound (A.66), it is easy to prove (A.62) with a fixed point argument. This proves the stability of (3.34)

**Lemma A.6.** *Suppose that  $\text{Im } w \sim 1$  and  $|m_c| \sim \text{Im } m_c \sim 1$ . Then  $|\partial_m f(\sqrt{w}, m_c)| \geq c$  for some constant  $c > 0$ .*

*Proof.* Using (2.13),  $m_c = \sqrt{w}(1 + m_{1c})$  and the conditions  $\text{Im } w \sim 1$ ,  $\text{Im } m_c \sim 1$ , we can get that

$$\left| \frac{\partial_{\sqrt{w}} f(\sqrt{w}, m_c)}{\partial_m f(\sqrt{w}, m_c)} \right| = \left| \frac{\partial m_c}{\partial \sqrt{w}} \right| \leq C \Rightarrow |\partial_{\sqrt{w}} f(\sqrt{w}, m_c)| \leq C |\partial_m f(\sqrt{w}, m_c)|, \quad (\text{A.69})$$

for some constant  $C > 0$ . Now we assume that  $|\partial_m f(\sqrt{w}, m_c)|$  can be arbitrarily small. Then  $|\partial_{\sqrt{w}} f(\sqrt{w}, m_c)|$  can also be arbitrarily small. Denote  $a := \partial_m f(\sqrt{w}, m_c)$  and  $b := \partial_{\sqrt{w}} f(\sqrt{w}, m_c)$ . Using (A.23) and (A.31), we get that

$$a = \frac{\sqrt{w}}{m_c} - \frac{m_c}{N} \sum_{i=1}^n l_i s_i \frac{\sqrt{w} (m_c^2 - |z|^2)^2 + 2s_i |z|^2 m_c}{[-(s_i + |z|^2)m_c^2 + |z|^4 + \sqrt{w} (m_c^3 - |z|^2 m_c)]^2} \quad (\text{A.70})$$

and

$$b = -1 - \frac{m_c^2}{N} \sum_{i=1}^n l_i s_i \frac{(m_c^2 - |z|^2)^2}{[-(s_i + |z|^2)m_c^2 + |z|^4 + \sqrt{w} (m_c^3 - |z|^2 m_c)]^2}. \quad (\text{A.71})$$

Using (A.70) and (A.71), we can get that

$$\frac{(\sqrt{w}m_c - |z|^2)|z|^2}{m_c} b - \frac{1}{2}(m_c^2 - |z|^2)(m_c a - \sqrt{w}b) = \frac{(|z|^2 - \sqrt{w}m_c)(m_c^2 + |z|^2)}{m_c}, \quad (\text{A.72})$$

where we use the equation  $f(\sqrt{w}, m_c) = 0$  in the derivation. By our assumption, the left-hand side of (A.72) can be arbitrarily small. For the right-hand side of (A.72), we have  $|m_c| \sim 1$  and  $|\sqrt{w}m_c - |z|^2| \sim 1$  (because  $\text{Im}(\sqrt{w}m_c) = \text{Im}(w + wm_{1c}) \sim 1$ ). Thus if  $|m_c - i|z| \geq c'$  for some constant  $c' > 0$ , we have  $|m^2 + |z|^2| \sim 1$ , and

$$\left| \frac{(\sqrt{w}m_c - |z|^2)|z|^2}{m_c} b - \frac{1}{2}(m_c^2 - |z|^2)(m_c a - \sqrt{w}b) \right| \sim 1,$$

which gives a contradiction. Thus we must have a lower bound  $|\partial_m f(\sqrt{w}, m_c)| \geq c$  if  $|m - i|z| \geq c'$ .

We still need to deal with the case where  $|m_c - i|z| \leq c'$  for some sufficiently small  $c'$ . Notice  $|z| \sim 1$  in this case. Then we have

$$\frac{\partial f}{\partial \sqrt{w}}(\sqrt{w}, i|z|) = -1 + \frac{|z|^2}{N} \sum_{i=1}^n l_i s_i \frac{4|z|^4}{[(s_i + |z|^2)|z|^2 + |z|^4 - 2i\sqrt{w}|z|^3]^2}. \quad (\text{A.73})$$

Denote  $L_i := (s_i + |z|^2)|z|^2 + |z|^4 - 2i\sqrt{w}|z|^3$ . Since  $i\sqrt{w} = i(x + iy) = ix - y$  with  $x, y > 0$  and  $x, y \sim 1$ , we have  $\text{Re } L_i > 0$ ,  $\text{Im } L_i < 0$  and  $|\text{Re } L_i|, |\text{Im } L_i| \sim 1$ . Furthermore,  $\text{Im } L_i^2 < 0$  and  $|\text{Im } L_i^2| \sim 1$ . Thus each fraction  $4|z|^4/L_i^2$  in (A.73) has positive imaginary part and all the imaginary parts have order 1. Therefore

$$\left| \frac{\partial f}{\partial \sqrt{w}}(\sqrt{w}, i|z|) \right| \geq \text{Im} \left[ \frac{\partial f}{\partial \sqrt{w}}(\sqrt{w}, i|z|) \right] \sim 1.$$

Then by (A.69), we get that  $|\partial_m f(\sqrt{w}, i|z|)| \geq c$  for some  $c > 0$ . Using (3.29), it is easy to see that

$$\partial_m f(\sqrt{w}, m_c) = \partial_m f(\sqrt{w}, i|z|) + O(|m_c - i|z|).$$

Thus in the case  $|m_c - i|z| \rightarrow 0$ , we still can find  $c > 0$  such that  $|\partial_m f(\sqrt{w}, m_c)| \geq c$ .  $\square$

**Lemma A.7.** Suppose that  $w \in \mathbf{D}_k^b$  and  $\text{Im } w \leq \epsilon$ . Then for sufficiently small  $\epsilon > 0$ , we have  $|\partial_m f(\sqrt{w}, m_c)| \sim 1$ .

*Proof.* By (3.22) and (3.29), if  $|w| \sim 1$  and  $\text{Im } m \sim 1$ , we have  $\partial_{\sqrt{w}} \partial_m f(w, m_c) = O(1)$  and  $\partial_m^2 f(w, m_c) = O(1)$ . Denote  $w = E + i\eta$ . Taking the imaginary part of the following equation

$$0 = f(\sqrt{E}, m_c(E)) = -\sqrt{E} + m_c + E^{-1/2} + \frac{1}{N} \sum_{i=1}^n l_i s_i \left( \frac{A_i}{m_c - a_i} + \frac{B_i}{m_c - b_i} + \frac{C_i}{m_c + c_i} \right), \quad (\text{A.74})$$



and noticing that  $A_i, B_i, C_i$  and  $a_i, b_i, c_i$  are all positive real numbers for real  $E$ , we get

$$\frac{1}{N} \sum_{i=1}^n l_i s_i \left( \frac{A_i}{|m_c - a_i|^2} + \frac{B_i}{|m_c - b_i|^2} + \frac{C_i}{|m_c + c_i|^2} \right) = 1. \quad (\text{A.75})$$

Using the above equation, we get

$$\begin{aligned} \partial_m f(\sqrt{E}, m_c(E)) &= 1 - \frac{1}{N} \sum_{i=1}^n l_i s_i \left[ \frac{A_i}{(m_c - a_i)^2} + \frac{B_i}{(m_c - b_i)^2} + \frac{C_i}{(m_c + c_i)^2} \right] \\ &= \frac{1}{N} \sum_{i=1}^n l_i s_i \left[ \frac{A_i}{|m_c - a_i|^2} - \frac{A_i}{(m_c - a_i)^2} + \frac{B_i}{|m_c - b_i|^2} - \frac{B_i}{(m_c - b_i)^2} + \frac{C_i}{|m_c + c_i|^2} - \frac{C_i}{(m_c + c_i)^2} \right]. \end{aligned} \quad (\text{A.76})$$

We look at, for example, the term

$$\frac{A_i}{|m_c - a_i|^2} - \frac{A_i}{(m_c - a_i)^2} = \frac{A_i}{|m_c - a_i|^2} (1 - e^{-2i\theta_i}),$$

where  $m_c - a_i := |m_c - a_i|e^{i\theta_i}$ . Using  $\text{Im } m_c \sim 1$ , it is easy to see that  $\text{Re}(1 - e^{-2i\theta_i}) \geq c'$  for some constant  $c' > 0$ . Applying the same estimates to the  $B, C$  terms in (A.76), we get

$$\left| \partial_m f(\sqrt{E}, m_c(E)) \right| \geq \text{Re} \left[ \partial_m f(\sqrt{E}, m_c(E)) \right] \geq c \quad (\text{A.77})$$

for some constant  $c > 0$ .

Now for  $w = E + i\eta$  with  $\eta \leq \epsilon$ , we can expand  $\partial_m f(\sqrt{w}, m_c(w))$  around  $\partial_m f(\sqrt{E}, m_c(E))$ ,

$$\partial_m f(\sqrt{w}, m_c(w)) = \partial_m f(E, m_c(E)) + O(\eta),$$

where we use (3.29). Combing with (A.77), we see that  $|\partial_m f(w, m_c(w))| \sim 1$  for small enough  $\epsilon$ .  $\square$

*Case 2:* We mimic the argument in the proof of Case 1. We see that it suffices to prove  $|\alpha| + |\partial_u \alpha| \leq C$  and  $|\beta| \sim 1$  for  $\alpha, \beta$  defined in (A.64) and (A.65) and  $|u - m_c| \leq (\log N)^{-1/3}$ . Using (3.29), it is not hard to prove that  $|\alpha| + |\partial_u \alpha| + |\beta| \leq C$ . What remains is the proof of the lower bound  $|\beta| \geq c$ . For the case  $\text{Im } w \sim 1$ , it follows from Lemma A.6. If  $w \rightarrow 0$  in the case  $|z|^2 \geq 1 + \tau$ , then  $m_c(w) = O(\sqrt{w}) \rightarrow 0$  by (3.23). Thus we can use (A.23) to get directly that

$$\partial_m f(\sqrt{w}, m_c) = 1 - |z|^{-2} + O(\sqrt{w}) \geq c.$$

Finally, we are left with the case  $E = \text{Re } w \sim 1$  and  $\eta = \text{Im } w \rightarrow 0$ . Using (2.13),  $m_c = \sqrt{w}(1 + m_{1c})$ ,  $|w| \sim 1$  and  $\text{dist}(E, \text{supp } \rho_{1c}) \geq \tau'$ , we can get that

$$\left| \frac{\partial_{\sqrt{w}} f(\sqrt{w}, m_c)}{\partial_m f(\sqrt{w}, m_c)} \right| = \left| \frac{\partial m_c}{\partial \sqrt{w}} \right| \leq C$$

for some constant  $C > 0$ . Thus it suffices to prove that  $|\partial_{\sqrt{w}} f(\sqrt{w}, m_c)|$  has a lower bound. Using (A.31) and noticing that  $m_c(E) \in \mathbb{R}$ , we get

$$\partial_{\sqrt{w}} f(\sqrt{E}, m_c(E)) = -1 - \frac{m_c^2}{N} \sum_{i=1}^n l_i s_i \frac{(m_c^2 - |z|^2)^2}{[-(s_i + |z|^2)m_c^2 + |z|^4 + \sqrt{E}(m_c^3 - |z|^2 m_c)]^2} \leq -1.$$

Expanding  $\partial_{\sqrt{w}}f(\sqrt{w}, m_c(w))$  around  $\partial_{\sqrt{w}}f(\sqrt{E}, m_c(E))$ , using (3.29) and  $|m_c(E+i\eta) - m_c(E)| \sim \eta$ , we get for  $\eta$  small

$$|\partial_{\sqrt{w}}f(\sqrt{w}, m_c)| \geq 1 + O(\eta) \geq c.$$

*Case 3:* The case  $\text{Im } w \geq \tau'$  can be proved with the same method as in the proof of case 1. Hence we only consider the case  $|w - e_k| \leq 2\tau'$  in the following. Note that  $|w| \sim 1$  in this case. Suppose

$$|w - e_k| \leq 2\tau', \quad |u - m_c| \leq (\log N)^{-1/3}. \quad (\text{A.78})$$

Then we claim that

$$|\alpha| \sim 1, \quad |\beta| \sim \sqrt{\kappa + \eta} \quad (\text{A.79})$$

for small enough  $\tau'$ . Using (A.78), (3.29), (2.19) and Lemma A.5, we can get that

$$\alpha = \frac{1}{2} \partial_m^2 f(\sqrt{e_k}, m_c(e_k)) + O(|w - e_k|^{1/2} + (\log N)^{-1/3}) \sim 1.$$

To prove the estimate for  $\beta$ , we use (2.17), (3.29) and Lemma A.5, to get

$$\begin{aligned} \beta &= \int_{e_k}^w \frac{d}{dw'} \partial_m f(\sqrt{w'}, m_c(w')) dw' = \int_{e_k}^w \frac{\partial_{\sqrt{w'}} \partial_m f(\sqrt{w'}, m_c(w'))}{2\sqrt{w'}} dw' + \int_{e_k}^w \partial_m^2 f(\sqrt{w'}, m_c(w')) \frac{dm_c(w')}{dw'} dw' \\ &= \int_{e_k}^w \frac{\partial_{\sqrt{w'}} \partial_m f(\sqrt{e_k}, m_c(e_k)) + O(|w - e_k|^{1/2})}{2\sqrt{w'}} dw' + \int_{m_c(e_k)}^{m_c(w)} \left[ \partial_m^2 f(\sqrt{e_k}, m_c(e_k)) + O(|w - e_k|^{1/2}) \right] dm \\ &= \partial_m^2 f(\sqrt{e_k}, m_c(e_k)) (m_c(w) - m_c(e_k)) + O(|w - e_k|). \end{aligned} \quad (\text{A.80})$$

Thus we conclude for small enough  $\tau'$  that

$$|\beta| \sim |w - e_k|^{1/2} \sim \sqrt{\kappa + \eta}.$$

With the estimate (A.79), we now proceed exactly as in the proof of [4, Lemma 4.5], by solving the quadratic equation (A.63) for  $u - m_c$  explicitly. We select the correct solution by a continuity argument using that (A.62) holds by assumption at  $z + iN^{-10}$ . The second assumption of (A.78) is obtained by continuity from the estimate on  $|u - m_c|$  at the neighboring point  $z + iN^{-10}$ . We refer to [4, Lemma 4.5] for the full details. This concludes the proof.

*Case 4:* The case when  $\text{Im } w \geq \tau'$  can be proved using the same method as in the proof of Case 1. Now we are left with the case  $|w| \leq 2\tau'$  for some sufficiently small  $\tau'$ . First we assume  $|z| \geq c > 0$  for some small  $c > 0$ . Then mimicking the argument in the proof of Case 1, we see that it suffices to prove  $|\alpha| + |\partial_u \alpha| \leq C$  and  $|\beta| \sim 1$  when  $|u - m_c| \leq (\log N)^{-1/3}$ . Using (3.29), it is not hard to prove that  $|\alpha| + |\partial_u \alpha| + |\beta| \leq C$ . The lower bound  $|\beta| \geq c$  can be obtained easily from (A.32).

Then suppose  $|z|^2 < c$ , but  $|w|^{1/2} + |z|^2 \geq \epsilon$ . According to (A.38) and using that  $|i|z|^2 + \sqrt{wt_0}| \sim |w|^{1/2} + |z|^2$ , we can verify that

$$\beta = \partial_m f(\sqrt{w}, m_c(w)) \sim |w|^{1/2} + |z|^2 \sim 1.$$

It is also easy to verify that

$$\partial_m^2 f(\sqrt{w}, m_c(w)) = O(|w|^{1/2} + |z|^2), \quad \partial_m^3 f(\sqrt{w}, m_c(w)) = O(|w|^{1/2} + |z|^2).$$

Hence if  $|u - m_c| \leq (\log N)^{-1/3}$ , we have

$$\alpha = \frac{1}{2} \partial_m^2 f(\sqrt{w}, m_c(w)) + O\left(\partial_m^3 f(\sqrt{w}, m_c(w)) (\log N)^{-1/3}\right) = O(|w|^{1/2} + |z|^2).$$

With a fixed point argument, we get (A.62).

*Case 5:* Again we following the arguments in the proof of Case 1. However, instead of  $f(\sqrt{w}, m)$ , we shall study  $\Upsilon(w, m_1)$  in (3.35) directly. We take over the notations in Definition 3.9 and abbreviate  $R := \Upsilon(w, u_1)$ , so that  $|R| \leq \delta$ . Then we write the equation  $\Upsilon(w, u_1) - \Upsilon(w, m_{1c}) = R$  as

$$\alpha(u_1)(u_1 - m_{1c})^2 + \beta(u_1 - m_{1c}) = R, \quad (\text{A.81})$$

where we use the same symbol as in (A.63) for notational convenience. As in Case 1, we have  $\beta = \partial_{m_1} \Upsilon(w, m_{1c})$ , and we can evaluate that  $|\alpha| + |\partial_{u_1} \alpha| \leq C$  for  $w \in \mathbf{D}_L$  and  $u_1$  satisfying  $|u_1 - m_{1c}| \ll |m_{1c}|$ . Now to conclude (3.39), it suffices to prove  $|\beta| \sim 1$  for  $w \in \mathbf{D}_L$ . In fact using (3.35), we obtain that

$$\beta = 1 + O(\eta^{-1}) \sim 1,$$

for  $\eta \geq \zeta^{-1}$ . This concludes the proof.

*Proof of Lemma 2.2.* The fact that  $\rho_{1c}$  has compact support follows from Lemma 2.3;  $\rho_{1c}$  being integrable follows from Lemma A.4. Note that in proving Lemmas 2.3 and A.4, we do not make the regularity assumptions in Definition 2.4. It remains to show that for fixed  $w \in \mathbb{C}_+$  and  $|z| \neq 1$ , there exists a unique  $m_{1c}(w) \in \mathbb{C}_+$  satisfying equation (2.11). This follows from the  $\eta \sim 1$  case in the proof of *Case 1* in this section.  $\square$

*Remark:* The estimate (3.29) has been used repeatedly during the proof of Lemma 3.10. Here we remark that it also gives the stability of the regularity conditions in Definition 2.4 under perturbations of  $|z|$  and  $\rho_\Sigma$ . For example, we define the shifted empirical spectral density

$$\rho_{\Sigma, t} := \frac{1}{N \wedge M} \sum_{i=1}^{N \wedge M} \delta_{\sigma_i + t}, \quad (\text{A.82})$$

and the associated  $m_c(w, t)$  and function  $f(\sqrt{w}, m, t)$ . Given a regular edge  $e_k$ , it satisfies that

$$f(\sqrt{e_k}, m_k, t = 0) = 0, \quad \partial_m f(\sqrt{e_k}, m_k, t = 0) = 0.$$

where we denote  $m_k := m_c(e_k)$ . We have the Jacobian

$$J := \det \begin{pmatrix} \partial_{\sqrt{w}} f & \partial_m f \\ \partial_{\sqrt{w}} \partial_m f & \partial_m^2 f \end{pmatrix}_{(\sqrt{w}, m, t) = (\sqrt{e_k}, m_k, 0)} = \partial_{\sqrt{w}} f(\sqrt{e_k}, m_k, 0) \partial_m^2 f(\sqrt{e_k}, m_k, 0).$$

By (A.31), we have  $|\partial_{\sqrt{w}} f(\sqrt{e_k}, m_k, 0)| \geq 1$ . Combining with (2.19), we get  $|J| \geq \epsilon$ . Using (3.29), we can verify that  $\partial_t f(\sqrt{e_k}, m_k, 0) = O(1)$  and  $\partial_t \partial_m f(\sqrt{e_k}, m_k, 0) = O(1)$ . Thus if we regard  $e_k$  and  $m_k$  as functions of  $t$ , then  $\partial_t m_k(t = 0) = O(1)$  and  $\partial_t e_k(t = 0) = O(1)$  by the implicit function theorem. Then it is easy to verify

$$\begin{aligned} \partial_m^2 f(\sqrt{e_k(t)}, m_c(e_k, t)) &= \partial_m^2 f(\sqrt{e_k}, m_c(e_k)) + O(t), \\ |m_c(e_k, t) - a_i(e_k, t)| &= |m_c(e_k) - a_i(e_k)| + O(t), \end{aligned}$$

and the similar estimates for  $|m_c - b_i|$  and  $|m_c + c_i|$ . Thus if Definition 2.4 (i) holds for some  $\rho_\Sigma$ , then it holds for all  $\rho_{\Sigma, t}$  provided that  $t$  is small enough.

Now given a regular bulk component  $[e_{2k}, e_{2k-1}]$  and  $E \in [e_{2k} + \tau', e_{2k-1} - \tau']$ . Differentiating the equation  $f(\sqrt{E}, m_c(E, t), t) = 0$  in  $t$  yields

$$\partial_t m_c(E, t) = - \frac{\partial_t f(\sqrt{E}, m_c(E, t), t)}{\partial_m f(\sqrt{E}, m_c(E, t), t)}.$$

By (3.29), we find that  $\partial_t f(\sqrt{E}, m_c(E), 0) = O(1)$ , while by (A.66),  $|\partial_m f(\sqrt{E}, m_c(E), 0)| = \beta \sim 1$ . Thus  $\partial_t m_c(E, 0) = O(1)$ . A simple extension of this argument shows that  $m_c(E, t) = m_c(E) + O(t)$  and hence  $\text{Im } m_c(E, t)$  is bounded from below by some  $c' = c'(\tau, \tau')$ . Thus we conclude that if Definition 2.4 (ii) holds for some  $\rho_\Sigma$ , then it holds for all  $\rho_{\Sigma, t}$  with  $t$  in some fixed small interval around zero. Obviously, the above arguments also work for the perturbation of  $|z|$ .

## B Proof of Lemma 4.9

Our proof of (4.59) is an extension of [4, Lemma 4.9], [7, Lemma 7.3] and [14, Theorem 4.7]. Here we only prove the bound for  $\|Z\|$ . The proof for  $\|\langle Z \rangle\|$  is exactly the same. For  $i \in \mathcal{I}_1$ , we define  $P_i = \mathbb{E}_{[i]}$  and  $Q_i = 1 - P_i$ . Recall that  $Z_{[i]} = Q_i G_{[ii]}^{-1}$ , we need to prove that

$$[Z] = \frac{1}{N} \sum_{i=1}^N \pi_{[i]} \left( Q_i G_{[ii]}^{-1} \right) \pi_{[i]} < |w|^{-1/2} \Phi_o^2,$$

for  $w \in \mathbf{D}$ . For  $J \subset \mathcal{I}$ , we define  $\pi_{[i]}^{[J]}$  by replacing  $m_{1,2}$  in (2.36) with  $m_{1,2}^{[J]}$  defined in (4.6). As in (4.58), we can prove that  $|m_{1,2}^{[i]} - m_{1,2}| < |w|^{-1/2} \Phi_o^2$ , which further gives that

$$[Z] = \frac{1}{N} \sum_{i=1}^N \pi_{[i]}^{[i]} \left( Q_i G_{[ii]}^{-1} \right) \pi_{[i]}^{[i]} + O_{<} \left( |w|^{-1/2} \Phi_o^2 \right) = \frac{1}{N} \sum_{i=1}^N Q_i \left( \pi_{[i]}^{[i]} G_{[ii]}^{-1} \pi_{[i]}^{[i]} \right) + O_{<} \left( |w|^{-1/2} \Phi_o^2 \right).$$

Thus if we abbreviate  $B_i := |w|^{1/2} Q_i \left( \pi_{[i]}^{[i]} G_{[ii]}^{-1} \pi_{[i]}^{[i]} \right)$ , it suffices to prove that  $B := N^{-1} \sum_i B_i < \Phi_o^2$ . We estimate  $B$  by bounding the  $p$ -th moment of its norm by  $\Phi_o^{2p}$  for  $p = 2n$  with  $n \in \mathbb{N}$ , i.e.  $\mathbb{E} \|B\|^p < \Phi_o^{2p}$ . The lemma then follows from the Chebyshev's inequality. Using  $\|KK^\dagger\| = \|K\|^2$  for any square matrix  $K$ , we get that for  $p = 2n$ ,

$$\text{Tr}(BB^\dagger)^n \geq \|BB^\dagger\|^n = \|B\|^{2n}.$$

Thus it suffices to prove that

$$\mathbb{E} \text{Tr}(BB^\dagger)^{p/2} < \Phi_o^{2p}, \quad \text{for } p = 2n. \quad (\text{B.1})$$

This estimate can be proved with the same method in [14, Appendix B], with the only complication being that  $\pi_{[i]}$  is random and depends on  $i$ . In principle, this can be handle by using (3.9) and (3.10) to put any indices  $j, k, \dots \in \mathcal{I}_1$  (that we wish to include) into the superscripts of  $\pi_{[i]}$ . This leads to a minor modification of the proof in [14, Appendix B]. Here we describe the basic ideas of the proof, without writing down all the details.

The proof is based on a decomposition of the space of random variables using  $P_s$  and  $Q_s$ . It is evident that  $P_s$  and  $Q_s$  are projections,  $P_s + Q_s = 1$  and all of these projections commute with each other. For a set  $J \subset \mathcal{I}$ , we denote  $P_J := \prod_{s \in J} P_s$  and  $Q_J := \prod_{s \in J} Q_s$ . Let  $p = 2n$  and introduce the shorthand notation  $\tilde{B}_{k_s} := B_{k_s}$  for  $s \leq p$  odd and  $\tilde{B}_{k_s} := B_{k_s}^\dagger$  for  $s \leq p$  even. Then we get

$$\mathbb{E} \text{Tr}(BB^\dagger)^{p/2} = \frac{1}{N^p} \sum_{k_1, k_2, \dots, k_p} \mathbb{E} \text{Tr} \prod_{s=1}^p \tilde{B}_{k_s} = \frac{1}{N^p} \sum_{k_1, k_2, \dots, k_p} \mathbb{E} \text{Tr} \prod_{s=1}^p \left( \prod_{r=1}^p (P_{k_r} + Q_{k_r}) \tilde{B}_{k_s} \right). \quad (\text{B.2})$$

Introducing the notations  $\mathbf{k} = (k_1, k_2, \dots, k_p)$  and  $\{\mathbf{k}\} = \{k_1, k_2, \dots, k_p\}$ , we can write

$$\mathbb{E}\text{Tr}(BB^\dagger)^{p/2} = \frac{1}{N^p} \sum_{\mathbf{k}} \sum_{I_1, \dots, I_p \subset \{\mathbf{k}\}} \mathbb{E}\text{Tr} \prod_{s=1}^p \left( P_{I_s^c} Q_{I_s} \tilde{B}_{k_s} \right). \quad (\text{B.3})$$

Following [14, Appendix B], we claim that to conclude (B.1) it suffices to prove that for  $k \in I$

$$\|Q_I B_k\| < \Phi_o^{|I|}. \quad (\text{B.4})$$

As in [14, Appendix B], it is not hard to prove for  $k \in I$ ,

$$|w|^{-1/2} \|Q_I G_{[kk]}^{-1}\| < \Phi_o^{|I|}. \quad (\text{B.5})$$

Now we extend the proof to obtain the estimate (B.4). For the case  $|I| = 1$  (i.e.  $I = \{k\}$ ),

$$\|B_k\| = |w|^{1/2} \|\pi_{[i]}^{[i]} Z_{[k]} \pi_{[i]}^{[i]}\| \leq |w|^{-1/2} \|Z_{[k]}\| < \Phi_o,$$

where we can prove  $\|Z_{[k]}\| < |w|^{1/2} \Phi_o$  by modifying the proof in Lemma 4.4. For the case  $|I| \geq 2$ , WLOG, we may assume  $k = 1$  and  $I = \{1, \dots, t\}$  with  $t \geq 2$ . It is enough to prove that

$$|w|^{1/2} \|Q_t \dots Q_2 \pi_{[1]}^{[1]} G_{[11]}^{-1} \pi_{[1]}^{[1]}\| < \Phi_o^t. \quad (\text{B.6})$$

We take the  $t = 3$  as an example to describe the ideas for the proof of (B.6). Using (3.9), we get

$$\pi_{[1]}^{[1]} = \pi_{[1]}^{[12]} + |w|^{1/2} \epsilon_{11}^{[1]} \pi_{[1]}^{[12]} A_1 \pi_{[1]}^{[12]} + |w|^{1/2} \epsilon_{11}^{[1]} \pi_{[1]}^{[12]} A_2 \pi_{[1]}^{[12]} + \text{error}_{1,2}, \quad (\text{B.7})$$

where  $\epsilon_{11}^{[1]}$  and  $\epsilon_{11}^{[1]}$  are entries of

$$\epsilon_{[1]}^{[1]} := |w|^{1/2} \left( \frac{G_{[22]}^{[1]}}{N} + \frac{1}{N} \sum_{k \notin \{1,2\}} G_{[k2]}^{[1]} \left( G_{[22]}^{[1]} \right)^{-1} G_{[2k]}^{[1]} \right) < \Phi_o^2,$$

$A_{1,2}$  are deterministic matrices with operator norm  $O(1)$ , and  $\|\text{error}_{1,2}\| < |w|^{-1/2} \Phi_o^4$ . Then we get

$$\begin{aligned} \pi_{[1]}^{[1]} G_{[11]}^{-1} \pi_{[1]}^{[1]} &= \pi_{[1]}^{[12]} G_{[11]}^{-1} \pi_{[1]}^{[12]} + |w|^{1/2} \epsilon_{11}^{[1]} \pi_{[1]}^{[12]} A_1 \pi_{[1]}^{[12]} G_{[11]}^{-1} \pi_{[1]}^{[12]} + |w|^{1/2} \epsilon_{11}^{[1]} \pi_{[1]}^{[12]} A_2 \pi_{[1]}^{[12]} G_{[11]}^{-1} \pi_{[1]}^{[12]} \\ &\quad + |w|^{1/2} \pi_{[1]}^{[12]} G_{[11]}^{-1} \epsilon_{11}^{[1]} \pi_{[1]}^{[12]} A_1 \pi_{[1]}^{[12]} + |w|^{1/2} \pi_{[1]}^{[12]} G_{[11]}^{-1} \epsilon_{11}^{[1]} \pi_{[1]}^{[12]} A_2 \pi_{[1]}^{[12]} + O_{<}(|w|^{-1/2} \Phi_o^4). \end{aligned} \quad (\text{B.8})$$

We first handle the  $\pi_{[1]}^{[12]} G_{[11]}^{-1} \pi_{[1]}^{[12]}$  term. By (B.5)

$$Q_2 \pi_{[1]}^{[12]} G_{[11]}^{-1} \pi_{[1]}^{[12]} = \pi_{[1]}^{[12]} \left( Q_2 G_{[11]}^{-1} \right) \pi_{[1]}^{[12]} < |w|^{-1/2} \Phi_o^2.$$

For the remaining term, we first expand  $\pi_{[1]}^{[12]} = \pi_{[1]}^{[123]} + O_{<}(|w|^{-1/2} \Phi_o^2)$  and use (B.5) to get

$$Q_3 Q_2 \pi_{[1]}^{[12]} G_{[11]}^{-1} \pi_{[1]}^{[12]} = \pi_{[1]}^{[123]} \left( Q_3 Q_2 G_{[11]}^{-1} \right) \pi_{[1]}^{[123]} + O_{<}(|w|^{-1/2} \Phi_o^4) < |w|^{-1/2} \Phi_o^3.$$

Then we deal with the second terms in (B.8). We first expand  $\epsilon_{[1]}^{[1]} = e_{[1]}^{[3]} + O_{<}(\Phi_o^3)$ , where

$$e_{[1]}^{[3]} := |w|^{1/2} \left( \frac{G_{[22]}^{[13]}}{N} + \frac{1}{N} \sum_{k \notin \{1,2,3\}} G_{[k2]}^{[13]} \left( G_{[22]}^{[13]} \right)^{-1} G_{[2k]}^{[13]} \right).$$

Using the similar arguments as above, we have

$$\begin{aligned} Q_3|w|^{1/2}e_{11}^{[3]}\pi_{[1]}^{[12]}A_1\pi_{[1]}^{[12]}G_{[11]}^{-1}\pi_{[1]}^{[12]} &= |w|^{1/2}e_{11}^{[3]}\pi_{[1]}^{[123]}A_1\pi_{[1]}^{[123]}\left(Q_3G_{[11]}^{-1}\right)\pi_{[1]}^{[123]} + O_{<}(|w|^{-1/2}\Phi_o^4) \\ &< |w|^{-1/2}\Phi_o^4. \end{aligned}$$

Thus we have

$$Q_2Q_3|w|^{1/2}e_{11}^{[1]}\pi_{[1]}^{[12]}A_1\pi_{[1]}^{[12]}G_{[11]}^{-1}\pi_{[1]}^{[12]} < |w|^{-1/2}\Phi_o^3.$$

Obviously this estimate works for the rest of the terms in (B.8). This proves (B.6) when  $t = 3$ .

We can continue in this manner for a general  $t$ . At the  $l$ -th step, we expand the leading order terms using (3.9) and (3.10), and after applying  $Q_l \dots Q_3Q_2$  on them, the number of  $\Phi_o$  factors increases by one at each step by (B.5). Trough induction we can prove (B.6). In fact the expansions can be performed in a systematic way using the method in [14, Appendix B], and we leave the details to the reader. Also we remark that similar techniques are used in the proof in Section 5, and we choose to present the details there (in fact the proof here is much easier than the one in Section 5).

## References

- [1] Z. D. Bai. Circular law. *Ann. Probab.*, 25(1):494–529, 1997.
- [2] Z. D. Bai and J. W. Silverstein. *Spectral Analysis of Large Dimensional Random Matrices*, volume 2 of *Mathematics Monograph Series*. Science Press, Beijing, 2006.
- [3] Z. Bao, G. Pan, and W. Zhou. Universality for the largest eigenvalue of sample covariance matrices with general population. *Ann. Statist.*, 43(1):382–421, 2015.
- [4] A. Bloemendal, L. Erdős, A. Knowles, H.-T. Yau, and J. Yin. Isotropic local laws for sample covariance and generalized Wigner matrices. *Electron. J. Probab.*, 19(33):1–53, 2014.
- [5] A. Bloemendal, A. Knowles, H.-T. Yau, and J. Yin. On the principal components of sample covariance matrices. *Prob. Theor. Rel. Fields*, 164(1):459–552, 2016.
- [6] A. Borodin and C. D. Sinclair. The Ginibre ensemble of real random matrices and its scaling limits. *Commun. Math. Phys.*, 291(1):177–224, 2009.
- [7] P. Bourgade, H.-T. Yau, and J. Yin. Local circular law for random matrices. *Probab. Theory Relat. Fields*, 159:545–595, 2014.
- [8] P. Bourgade, H.-T. Yau, and J. Yin. The local circular law II: the edge case. *Probab. Theory Relat. Fields*, 159(3):619–660, 2014.
- [9] K. R. Davidson and S. J. Szarek. Local operator theory, random matrices and banach spaces. volume 1 of *Handbook of the Geometry of Banach Spaces*, pages 317 – 366. North-Holland, Amsterdam, 2001.
- [10] A. Edelman. The probability that a random real gaussian matrix has  $k$  real eigenvalues, related distributions, and the circular law. *J. Multivar. Anal.*, 60(2):203 – 232, 1997.
- [11] N. El Karoui. Tracy-widom limit for the largest eigenvalue of a large class of complex sample covariance matrices. *Ann. Probab.*, 35(2):663–714, 2007.

- [12] L. Erdős, A. Knowles, and H.-T. Yau. Averaging fluctuations in resolvents of random band matrices. *Ann. Henri Poincaré*, 14:1837–1926, 2013.
- [13] L. Erdős, A. Knowles, H.-T. Yau, and J. Yin. Delocalization and diffusion profile for random band matrices. *Commun. Math. Phys.*, 323:367–416, 2013.
- [14] L. Erdős, A. Knowles, H.-T. Yau, and J. Yin. The local semicircle law for a general class of random matrices. *Electron. J. Probab.*, 18:1–58, 2013.
- [15] L. Erdős, A. Knowles, H.-T. Yau, and J. Yin. Spectral statistics of Erdős-Rényi graphs I: Local semicircle law. *Ann. Probab.*, 41(3B):2279–2375, 2013.
- [16] L. Erdős, B. Schlein, and H.-T. Yau. Local semicircle law and complete delocalization for Wigner random matrices. *Commun. Math. Phys.*, 287(2):641–655, 2008.
- [17] L. Erdős, H.-T. Yau, and J. Yin. Bulk universality for generalized Wigner matrices. *Probab. Theory Relat. Fields*, 154(1):341–407, 2012.
- [18] P. J. Forrester and T. Nagao. Eigenvalue statistics of the real Ginibre ensemble. *Phys. Rev. Lett.*, 99:050603, 2007.
- [19] J. Ginibre. Statistical ensembles of complex, quaternion, and real matrices. *J. Math. Phys.*, 6(3):440–449, 1965.
- [20] V. Girko. The circular law. *Russ. Teor. Veroyatnost. i Primenen.*, 29(4):669–679, 1984.
- [21] F. Götze and A. Tikhomirov. The circular law for random matrices. *Ann. Probab.*, 38(4):1444–1491, 2010.
- [22] A. Guionnet, M. Krishnapur, and O. Zeitouni. The single ring theorem. *Ann. Math.*, 174(2):1189–1217, 2011.
- [23] W. Hachem, A. Hardy, and J. Najim. Large complex correlated Wishart matrices: Fluctuations and asymptotic independence at the edges. *arXiv:1409.7548*.
- [24] A. Knowles and J. Yin. Anisotropic local laws for random matrices. *arXiv:1410.3516*.
- [25] A. Knowles and J. Yin. The isotropic semicircle law and deformation of Wigner matrices. *Comm. Pure Appl. Math.*, 66(11):1663–1749, 2013.
- [26] J. O. Lee and K. Schnelli. Tracy-widom distribution for the largest eigenvalue of real sample covariance matrices with general population. *arXiv:1409.4979*.
- [27] A. Litvak, A. Pajor, M. Rudelson, and N. Tomczak-Jaegermann. Smallest singular value of random matrices and geometry of random polytopes. *Adv. Math.*, 195(2):491 – 523, 2005.
- [28] M. L. Mehta. *Random matrices*, volume 142 of *Pure and Applied Mathematics*. Elsevier, Amsterdam, 3 edition, 2004.
- [29] A. Onatski. The tracy-widom limit for the largest eigenvalues of singular complex wishart matrices. *Ann. Appl. Probab.*, 18(2):470–490, 2008.
- [30] G. Pan and W. Zhou. Circular law, extreme singular values and potential theory. *J. Multivar. Anal.*, 101(3):645–656, 2010.

- [31] M. Rudelson. Invertibility of random matrices: norm of the inverse. *Ann. Math.*, 168(2):575–600, 2008.
- [32] M. Rudelson and R. Vershynin. The littlewood-offord problem and invertibility of random matrices. *Adv. Math.*, 218:600–633, 2008.
- [33] M. Rudelson and R. Vershynin. The smallest singular value of a random rectangular matrix. *Comm. Pure Appl. Math.*, 62:1707–1739, 2009.
- [34] M. Rudelson and R. Vershynin. Small ball probabilities for linear images of high-dimensional distributions. *Int. Math. Res. Notices*, 2015(19):9594–9617, 2015.
- [35] C. D. Sinclair. Averages over Ginibre’s ensemble of random real matrices. *Int. Math. Res. Not.*, 2007:1–15.
- [36] T. Tao and V. Vu. Random matrices: the circular law. *Commun. Contemp. Math.*, 10(2):261–307, 2008.
- [37] T. Tao and V. Vu. Random matrices: Universality of local spectral statistics of non-Hermitian matrices. *Ann. Probab.*, 43(2):782–874, 2015.
- [38] T. Tao, V. Vu, and M. Krishnapur. Random matrices: Universality of ESDs and the circular law. *Ann. Probab.*, 38(5):2023–2065, 2010.
- [39] J. Yin. The local circular law III: general case. *Probab. Theory Relat. Fields*, 160(3):679–732, 2014.