

# Supplemental Material

## Double Trouble in Double Descent: Bias and Variance(s) in the Lazy Regime

### Contents

<b>1</b>	<b>Further analytical results</b>	<b>2</b>
1.1	Asymptotic scalings . . . . .	2
1.2	Divide and Conquer approach . . . . .	2
1.3	Is it always better to be overparametrized ? . . . . .	3
<b>2</b>	<b>Statement of the Main Result</b>	<b>3</b>
2.1	Assumptions . . . . .	3
2.2	Results . . . . .	4
2.3	Explicit expression of $S^v, S^e, S^d$ . . . . .	5
2.4	Explicit expression of $P_{\Psi_1}, P_{\Psi_2}^v, P_{\Psi_3}^v, P_{\Psi_2}^e, P_{\Psi_3}^e, P_{\Psi_2}^d$ . . . . .	5
<b>3</b>	<b>Replica Computation</b>	<b>6</b>
3.1	Toolkit . . . . .	6
3.1.1	Gaussian integrals . . . . .	6
3.1.2	Replica representation of an inverse matrix . . . . .	6
3.2	The Random Feature model . . . . .	7
3.2.1	With a single learner . . . . .	7
3.2.2	Ensembling over $K$ learners . . . . .	7
3.2.3	Equivalent Gaussian Covariate Model . . . . .	8
3.3	Computation of the <i>vanilla</i> terms . . . . .	8
3.3.1	Averaging over the dataset . . . . .	9
3.3.2	Averaging over the deterministic noise . . . . .	10
3.3.3	Averaging over the random feature vectors . . . . .	10
3.3.4	Expression of the action and the prefactor . . . . .	10
3.3.5	Expression of the action and the prefactor in terms of order parameters . . . . .	11
3.3.6	Saddle point equations . . . . .	12
3.3.7	Fluctuations around the saddle point . . . . .	12
3.3.8	Expression of the vanilla terms . . . . .	14
3.4	Computation of the <i>ensembling</i> terms . . . . .	14
3.4.1	Expression of the action and the prefactor . . . . .	14
3.4.2	Expression of the action and the prefactor in terms of order parameters . . . . .	15
3.4.3	Expression of the ensembling terms . . . . .	16
3.5	Computation of the <i>divide and conquer</i> term . . . . .	16

# 1 Further analytical results

## 1.1 Asymptotic scalings

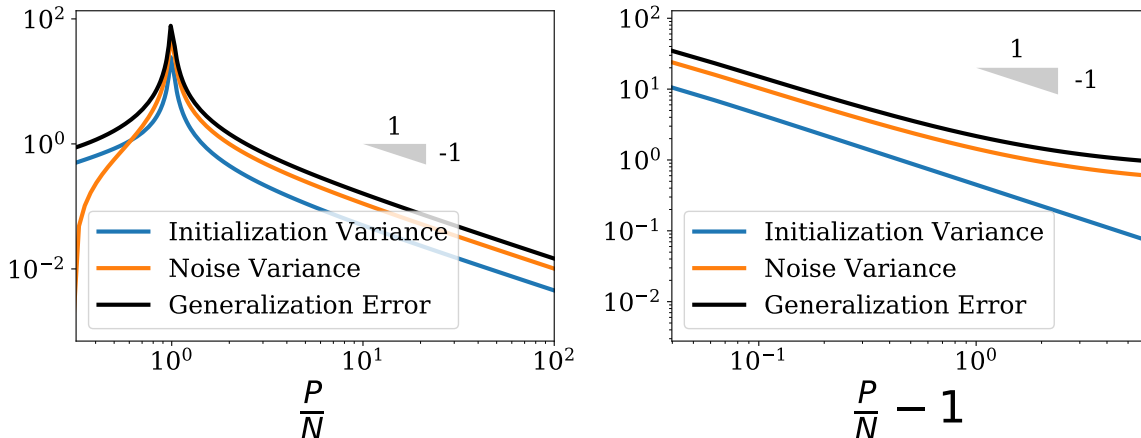


Figure 1: Log-log behaviour of the quantities of interest at **Left**:  $P/N \rightarrow \infty$  and **Right**:  $P \rightarrow N$  with  $\lambda = 10^{-5}$ ,  $N/D = 1$  and  $\tau = 1$ . In both cases, one observes an inverse scaling law.

Figure 1 (left) shows that the various terms entering the decomposition of the test error approach their asymptotic values at a rate  $(P/N)^{-1}$ . This scaling law is consistent with that found in [3] for real neural networks, where  $P$  is replaced by the width of the layers of the network. As for the divergence of the noise and initialization variances observed at the interpolation threshold, figure 1 (right) shows that they also follow an inverse power law  $(P/N - 1)^{-1}$  at vanishing regularization.

## 1.2 Divide and Conquer approach

As mentioned in the main text, another way to average the predictions of differently initialized learners is the *divide and conquer* approach [2]. In this framework, the data set is divided into  $K$  splits of size  $N/K$ . Each of the  $K$  differently initialized learner is trained on a distinct split. This approach is extremely useful for kernel learning [5], where the computational burden is in the inversion of the Gram matrix which is of size  $N \times N$ . In the random projection approach considered here, it does not offer any computational gain, however it is interesting how it affects the test error.

Within our framework, the test error can easily be calculated as:

$$\mathbb{E}_{\{\Theta^{(k)}\}, \mathbf{X}, \epsilon} [\mathcal{R}_{\text{RF}}] = F^2 (1 - 2\Psi_1^v) + \frac{1}{K} (F^2\Psi_2^v + \tau^2\Psi_3^v) + \left(1 - \frac{1}{K}\right) F^2\Psi_2^d, \quad (1)$$

where the effective number of data points which enters this formula is  $N_{\text{eff}} = N/K$  due to the splitting of the training set.

Comparing the previous expression with that obtained for ensembling (24) is instructive: here, increasing  $K$  replaces the *vanilla* terms  $\Psi_2^v, \Psi_3^v$  by the *divide and conquer* term  $\Psi_2^d$ . This shows that divide and conquer has a *denoising* effect: at  $K \rightarrow \infty$ , the effect of the additive noise on the labels is completely suppressed. This was not the case for ensembling. The price to pay is that  $N_{\text{eff}}$  decreases, hence one is shifted to the underparametrized regime.

In Figure 2, we see that the kernel limit error of the divide and conquer approach, i.e. the asymptotic value of the error at  $P/N \rightarrow \infty$ , is different from the usual kernel limit error, since the effective dataset is two times smaller at  $K = 2$ . The denoising effect of the divide and conquer approach is illustrated by the fact that its kernel limit error is higher at high SNR, but lower at low SNR. This is of practical relevance, and is much related to the beneficial effect of *bagging* in noisy dataset scenarios. The divide and conquer approach, which only differs from bagging by the fact that the different partitions of the dataset are disjoint, was shown to reach bagging-like performance in various setups such as decision trees and neural networks [1].

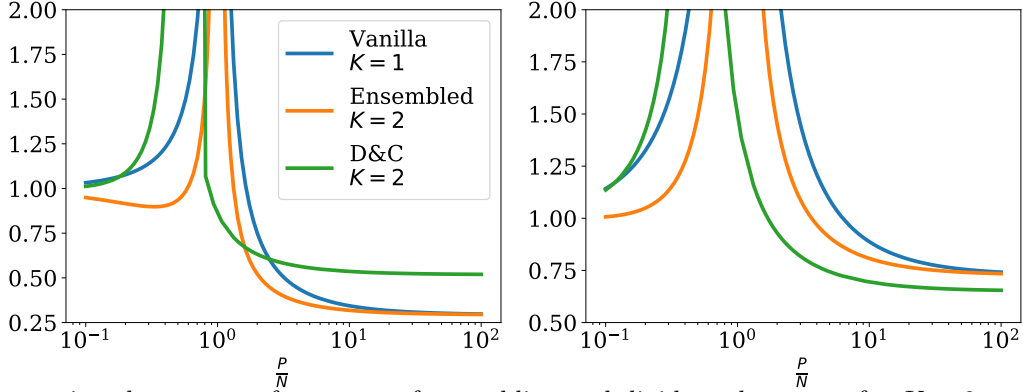


Figure 2: Comparison between performances of ensembling and divide and conquer for  $K = 2$  at different SNR. **Left:**  $\text{SNR} = \frac{F}{\tau} = 10$ . **Right:**  $\text{SNR} = \frac{F}{\tau} = 1$  Computations performed with fixed  $N/D = 1$ ,  $F = 1$  and  $\lambda = 10^{-5}$ .

### 1.3 Is it always better to be overparametrized ?

A common thought is that the double descent curve always reaches its minimum in the over-parametrized regime, leading to the idea that the corresponding model "cannot overfit". In this section, we show that this is not always the case. Three factors tend to shift the optimal generalization to the underparametrized regime: (i) increasing the numbers of learners from which we average the predictions,  $K$ , (ii) decreasing the signal-to-noise ratio (SNR),  $F/\tau$ , and (iii) decreasing the size of the dataset,  $N/D$ . In other words, when ensembling on a small, noisy dataset, one is better off using an underparametrized model.

These three effects are shown in figure 3. In the left panel, we see that as we increase  $K$ , the minimum of test error jumps to the underparametrized regime  $P < N$  for a high enough value of  $K$ . In the central/right panels, a similar effect occurs when decreasing the SNR or decreasing  $N/D$ .

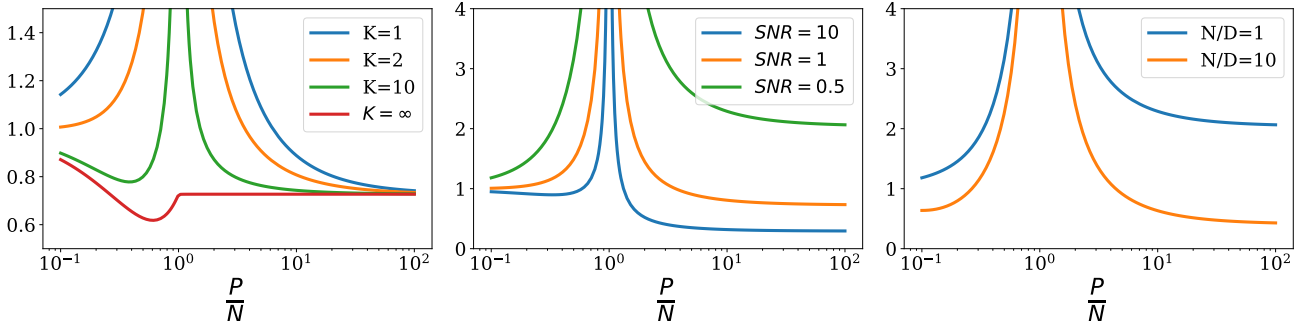


Figure 3: Generalisation error as a function of  $P/N$ : depending on the values of  $K$ ,  $F/\tau$  and  $N/D$ , optimal generalization can be reached in the underparametrized regime or the overparametrized regime. **Left:**  $F/\tau = 1$ ,  $N/D = 1$  and we vary  $K$ . This is the same as figure 5 in the main text, except that the higher noise causes the ensembling curve at  $K \rightarrow \infty$  to exhibit a *dip* in the underparametrized regime. **Center:**  $K = 2$ ,  $N/D = 1$  and we vary  $F/\tau$ . **Right:**  $F/\tau = 1$ ,  $K = 2$  and we vary  $N/D$ .

## 2 Statement of the Main Result

### 2.1 Assumptions

First, we state precisely the assumptions under which our main result is valid. Note, that these are the same as in [4].

**Assumption 1:**  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is a weakly differential function with derivative  $\sigma'$ . Assume there exists  $c_0, c_1 < \infty \in \mathbb{R}$  such that for all  $u \in \mathbb{R}$   $|\sigma(u)|, |\sigma'(u)| \leq c_0 e^{c_1|u|}$ . Then define:

$$\mu_0 = \mathbb{E}[\sigma(u)] \quad \mu_1 = \mathbb{E}[u\sigma(u)] \quad \mu_*^2 = \mathbb{E}[\sigma^2(u)] - \mu_0^2 - \mu_1^2, \quad (2)$$

where the expectation is over  $u \sim \mathcal{N}(0, 1)$ . To facilitate readability, we specialize to the case  $\mu_0 = 0$ . This simply amounts to a shift activation function  $\tilde{\sigma}$  of the network,  $\tilde{\sigma}(x) = \sigma(x) - \mu_0$ .

**Assumption 2:** We work in the high-dimensional limit, i.e. in the limit where the input dimension  $D$ , the hidden layer dimension  $P$  and the number of training points  $N$  go to infinity with their ratios fixed. That is:

$$N, P, D \rightarrow \infty, \quad \frac{P}{D} \equiv \psi_1 = \mathcal{O}(1), \quad \frac{N}{D} \equiv \psi_2 = \mathcal{O}(1). \quad (3)$$

This condition implies that, in the computation of the risk  $\mathcal{R}$ , we can neglect all the terms of order  $\mathcal{O}(1)$  in favour of the terms of order  $\mathcal{O}(D)$ .

**Assumption 3:** The labels are given by a linear ground truth, or teacher function:

$$y_\mu = f_d(\mathbf{X}_\mu) + \epsilon_\mu, \quad f_d(\mathbf{x}) = \langle \boldsymbol{\beta}, \mathbf{x} \rangle, \quad \|\boldsymbol{\beta}\| = F, \quad \epsilon_\mu \sim \mathcal{N}(0, \tau). \quad (4)$$

Note that as explained in [4], it is easy to add a non linear component to the teacher, but the latter would not be captured by the model (the student) in the regime  $N/D = \mathcal{O}(1)$ , and would simply amount to an extra noise term.

## 2.2 Results

Here we give the explicit form of the quantities appearing in our main result. In these expressions, the index  $a \in \{v, e, d\}$  distinguishes the *vanilla*, *ensembling* and *divide and conquer* terms.

$$\begin{aligned} \Psi_1 &= \frac{1}{D} \text{Tr} \left[ H [S^v]^{-1} H [P_{\Psi_1}] \right], \\ \Psi_2^v &= \frac{1}{D} \text{Tr} \left[ H [S^v]^{-1} H [P_{\Psi_2^v}] \right], \\ \Psi_3^v &= \frac{1}{D} \text{Tr} \left[ H [S^v]^{-1} H [P_{\Psi_3^v}] \right], \\ \Psi_2^e &= \frac{1}{D} \text{Tr} \left[ H [S^e]^{-1} H [P_{\Psi_2^e}] \right], \\ \Psi_3^e &= \frac{1}{D} \text{Tr} \left[ H [S^e]^{-1} H [P_{\Psi_3^e}] \right], \\ \Psi_2^d &= \frac{1}{D} \text{Tr} \left[ H [S^d]^{-1} H [P_{\Psi_2^d}] \right], \end{aligned}$$

where the Hessian matrix  $H[F]$ , for a given function  $F : (q, r, \tilde{q}, \tilde{r}) \mapsto \mathbb{R}$  is defined as:

$$H[F] = \begin{pmatrix} \frac{\partial F}{\partial q \partial q} & \frac{\partial F}{\partial q \partial r} & \frac{\partial F}{\partial q \partial \tilde{q}} & \frac{\partial F}{\partial q \partial \tilde{r}} \\ \frac{\partial F}{\partial q \partial r} & \frac{\partial F}{\partial r \partial r} & \frac{\partial F}{\partial r \partial \tilde{q}} & \frac{\partial F}{\partial r \partial \tilde{r}} \\ \frac{\partial F}{\partial q \partial \tilde{q}} & \frac{\partial F}{\partial r \partial \tilde{q}} & \frac{\partial F}{\partial \tilde{q} \partial \tilde{q}} & \frac{\partial F}{\partial \tilde{q} \partial \tilde{r}} \\ \frac{\partial F}{\partial q \partial \tilde{r}} & \frac{\partial F}{\partial r \partial \tilde{r}} & \frac{\partial F}{\partial \tilde{q} \partial \tilde{r}} & \frac{\partial F}{\partial \tilde{r} \partial \tilde{r}} \end{pmatrix} \Bigg|_{\substack{q=q^* \\ r=r^* \\ \tilde{r}=0 \\ \tilde{q}=0}},$$

with  $q^*$  and  $r^*$  being the solutions of the fixed point equation for the function  $S_0 : (q, r) \mapsto \mathbb{R}$  defined below:

$$\begin{cases} \frac{\partial S_0(q, r)}{\partial q} = 0 \\ \frac{\partial S_0(q, r)}{\partial r} = 0. \end{cases}$$

$$S_0(q, r) = \lambda \psi_1^2 \psi_2 q + \psi_2 \log \left( \frac{\mu_1^2 \psi_1 q}{\mu_1^2 \psi_1 r + 1} + 1 \right) + \frac{r}{q} + (1 - \psi_1) \log(q) + \psi_2 \log(\mu_1^2 \psi_1 r + 1) - \log(r).$$

The explicit expression of the above quantities in terms of  $(q, r, \tilde{q}, \tilde{r})$  is given below.

### 2.3 Explicit expression of $S^v, S^e, S^d$

Here we present the explicit formulas for  $S^v, S^e, S^d$ , which are defined as the functions  $(q, r, \tilde{q}, \tilde{r}) \mapsto \mathbb{R}$  such that:

$$\begin{aligned} S^v(q, r, \tilde{q}, \tilde{r}) &= 2(S_0(q, r) + \tilde{q}f^v(q, r) + \tilde{r}g^v(q, r)) \\ S^e(q, r, \tilde{q}, \tilde{r}) &= S_0(q, r) + \tilde{r}^2 f^e(q, r) + \tilde{q}^2 g^e(q, r) \\ S^d(q, r, \tilde{q}, \tilde{r}) &= S_0(q, r) + \tilde{r}^2 f^d(q, r) + \tilde{q}^2 g^d(q, r), \end{aligned}$$

where we defined the functions  $(q, r) \mapsto \mathbb{R}$ ,

$$\begin{aligned} f^v(q, r) &= \lambda\psi_1^2\psi_2 + \frac{\mu_\star^2\psi_1\psi_2}{\mu_\star^2\psi_1q + \mu_1^2\psi_1r + 1} + \frac{1 - \psi_1}{q} - \frac{r}{q^2}, \\ g^v(q, r) &= -\frac{\mu_\star^2\mu_1^2\psi_1^2\psi_2q}{(\mu_1^2\psi_1r + 1)(\mu_\star^2\psi_1q + \mu_1^2\psi_1r + 1)} + \frac{\mu_1^2\psi_1\psi_2}{\mu_1^2\psi_1r + 1} + \frac{1}{q} - \frac{1}{r}, \\ f^e(q, r) &= \frac{2r\mu_1^2\psi_1(1 + q\mu_\star^2\psi_1) + (1 + q\mu_\star^2\psi_1)^2 - r^2\mu_1^4\psi_1^2(-1 + \psi_2)}{r^2(1 + r\mu_1^2\psi_1 + q\mu_\star^2\psi_1)^2}, \\ g^e(q, r) &= \frac{\psi_1}{q^2}, \\ f^d(q, r) &= \frac{1}{r^2}, \\ g^d(q, r) &= \frac{\psi_1}{q^2}. \end{aligned}$$

### 2.4 Explicit expression of $P_{\Psi_1}, P_{\Psi_2}^v, P_{\Psi_3}^v, P_{\Psi_2}^e, P_{\Psi_3}^e, P_{\Psi_2}^d$

Here we present the explicit formulas for  $P_{\Psi_1}, P_{\Psi_2}^v, P_{\Psi_3}^v, P_{\Psi_2}^e, P_{\Psi_3}^e, P_{\Psi_2}^d$ , which are defined as the functions  $(q, r, \tilde{q}, \tilde{r}) \mapsto \mathbb{R}$  such that:

$$\begin{aligned} P_{\Psi_1} &= \psi_1\psi_2\mu_1^2 \left( M_X^{11} + \mu_1^2\mu_\star^2\psi_1^2 (M_X M_W^v M_X)^{11} + \mu_\star^2\psi_1 (M_X M_W^v)^{11} \right), \\ P_{\Psi_2}^v &= D\psi_1^2\psi_2(\mu_1^2\tilde{r} + \mu_\star^2\tilde{q}) \left[ \mu_1^2 P_{XX}^v - 2\mu_1^2\mu_\star^2\psi_1 P_{WX}^v + \mu_\star^2 P_{WW}^v \right], \\ P_{\Psi_3}^v &= D\psi_1^2\psi_2(\mu_1^2\tilde{r} + \mu_\star^2\tilde{q}) \left[ \mu_1^2 \left( M_X^{12} + \mu_1^2\mu_\star^2\psi_1^2 [M_X M_W^v M_X]^{12} \right) - 2\mu_1^2\mu_\star^2\psi_1 [M_X M_W^v]^{12} + \mu_\star^2 [M_W^v]^{12} \right], \\ P_{\Psi_2}^e &= D\psi_1^2\psi_2\mu_1^2\tilde{r} \left[ \mu_1^2 P_{XX}^e - 2\mu_1^2\mu_\star^2\psi_1 P_{WX}^e + \mu_\star^2 P_{WW}^e \right], \\ P_{\Psi_3}^e &= D\psi_1^2\psi_2\mu_1^2\tilde{r} \left[ \mu_1^2 \left( M_X^{12} + \mu_1^2\mu_\star^2\psi_1^2 [M_X M_W^e M_X]^{12} \right) - 2\mu_1^2\mu_\star^2\psi_1 [M_X M_W^e]^{12} + \mu_\star^2 [M_W^e]^{12} \right], \\ P_{\Psi_2}^d &= D\mu_1^2\psi_1\psi_2^2\tilde{r} \left[ \psi_1\mu_1^2 P_{XX} + 2\mu_\star^2\mu_1^2\psi_1^2 P_{WX} + \mu_\star^2\psi_1 P_{WW} \right], \end{aligned}$$

where we defined the scalars  $P_{XX}, P_{WX}, P_{WW}$  as follows:

$$\begin{aligned} P_{XX}^v &= \psi_2 N_X^{12} + M_X^{12} + 2\psi_2(\mu_1\mu_\star\psi_1)^2 [M_X N_X M_W^a]^{12} + (\mu_1\mu_\star\psi_1)^2 [M_X M_W^a M_X]^{12} \\ &\quad + \psi_2(\mu_1\mu_\star\psi_1)^4 [M_X M_W N_X M_W^a M_X]^{12}, \\ P_{WX}^v &= \psi_2 [N_X M_W^a]^{12} + [M_X M_W^a]^{12} + \psi_2(\mu_1\mu_\star\psi_1)^2 [M_X M_W^a N_X M_W^a]^{12}, \\ P_{WW}^v &= [M_W^a]^{12} + \psi_2(\mu_1\mu_\star\psi_1)^2 [M_W^a N_X M_W^a]^{12}, \\ P_{XX}^e &= P_{XX}^v, \\ P_{WX}^e &= P_{WX}^v, \\ P_{WW}^e &= P_{WW}^v, \\ P_{XX}^d &= \left( N_X^{d11} + 2(\mu_1\mu_\star\psi_1)^2 [N_X^d M_W^d M_X^d]^{11} + (\mu_1\mu_\star\psi_1)^4 [M_X^d M_W^d N_X^d M_W^d M_X^d]^{11} \right), \\ P_{WX}^d &= [N_X^d M_W^d]^{11} + (\mu_1\mu_\star\psi_1)^2 [M_X^d M_W^d N_X^d M_W^d]^{11}, \\ P_{WW}^d &= (\mu_1\mu_\star\psi_1)^2 [M_W^d N_X^d M_W^d]^{11}, \end{aligned}$$

and the  $2 \times 2$  matrices  $M_X, M_W, N_X$  as follows:

$$M_X^v = \begin{bmatrix} \frac{r}{1+\mu_1^2\psi_1 r} & \frac{\tilde{r}}{(1+\mu_1^2\psi_1 r)^2} \\ \frac{\tilde{r}}{(1+\mu_1^2\psi_1 r)^2} & \frac{r}{1+\mu_1^2\psi_1 r} \end{bmatrix}, \quad M_W^v = \begin{bmatrix} \frac{q(1+\mu_1^2\psi_1 r)}{1+2\mu_1^2\psi_1 r+\mu_*^2\psi_1 q} & \frac{q^2\mu_1^2\mu_*^2\psi_1^2\tilde{r}}{(1+\mu_1^2\psi_1 r+\mu_*^2\psi_1 q)^2} \\ \frac{q^2\mu_1^2\mu_*^2\psi_1^2\tilde{r}}{(1+\mu_1^2\psi_1 r+\mu_*^2\psi_1 q)^2} & \frac{q(1+\mu_1^2\psi_1 r)}{1+2\mu_1^2\psi_1 r+\mu_*^2\psi_1 q} \end{bmatrix}, \quad N_X^v = \frac{1}{(1+\mu_1^2\psi_1 r)^2} \begin{bmatrix} r & \tilde{r} \\ \tilde{r} & r \end{bmatrix},$$

$$M_X^e = M_X^v, \quad M_W^e = M_W^v + \frac{(1+r\mu_1^2\psi_1)^2\tilde{q}}{(1+\mu_1^2\psi_1 r+\mu_*^2\psi_1 q)^2} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad N_X^e = \frac{1}{(1+\mu_1^2\psi_1 r)^2} \begin{bmatrix} r & \tilde{r} \\ \tilde{r} & r \end{bmatrix},$$

$$M_X^d = \frac{r}{1+\mu_1^2\psi_1 r^2} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad M_W^d = \frac{q(1+\mu_1^2\psi_1 r)}{1+\mu_1^2\psi_1 r+\mu_*^2\psi_1 q} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad N_X^d = \frac{1}{(1+\mu_1^2\psi_1 r)^2} \begin{bmatrix} \tilde{r} & 0 \\ 0 & \tilde{r} \end{bmatrix}.$$

### 3 Replica Computation

#### 3.1 Toolkit

##### 3.1.1 Gaussian integrals

In order to obtain the main result for the generalisation error, we perform the averages over all the sources of randomness in the system in the following order: over the dataset  $X$ , then over the noise  $W$ , and finally over the random feature layers  $\Theta$ . Here are some useful formulaes used throughout the computations:

$$\begin{cases} \int e^{-\frac{1}{2}x_i G_{ij} x_j + J_i x_i} dx & = (\det G)^{-\frac{1}{2}} e^{\frac{1}{2}J_i G_{ij}^{-1} J_j}, \\ \int x_a e^{-\frac{1}{2}x_i G_{ij} x_j + J_i x_i} dx & = P_a^1 (\det G)^{-\frac{1}{2}} e^{\frac{1}{2}J_i G_{ij}^{-1} J_j}, \\ \int x_a x_b e^{-\frac{1}{2}x_i G_{ij} x_j + J_i x_i} dx & = P_{ab}^2 (\det G)^{-\frac{1}{2}} e^{\frac{1}{2}J_i G_{ij}^{-1} J_j}, \\ \int x_a x_b x_c e^{-\frac{1}{2}x_i G_{ij} x_j + J_i x_i} dx & = P_{abc}^3 (\det G)^{-\frac{1}{2}} e^{\frac{1}{2}J_i G_{ij}^{-1} J_j}, \\ \int x_a x_b x_c x_d e^{-\frac{1}{2}x_i G_{ij} x_j + J_i x_i} dx & = P_{abcd}^4 (\det G)^{-\frac{1}{2}} e^{\frac{1}{2}J_i G_{ij}^{-1} J_j}, \end{cases} \quad (5)$$

with

$$P_a^1 = [G^{-1}J]_a,$$

$$P_{ab}^2 = ((G^{-1})_{ab} + [G^{-1}J]_a [G^{-1}J]_b),$$

$$P_{abc}^3 = \sum_{\bar{a}, \bar{b}, \bar{c} \in \text{perm}(abc)} ((G^{-1})_{\bar{a}\bar{b}} [G^{-1}J]_{\bar{c}} + [G^{-1}J]_{\bar{a}} [G^{-1}J]_{\bar{b}} [G^{-1}J]_{\bar{c}}),$$

$$P_{abcd}^4 = \sum_{\bar{a}, \bar{b}, \bar{c}, \bar{d} \in \text{perm}(abcd)} \left( (G^{-1})_{\bar{a}\bar{b}} (G^{-1})_{\bar{c}\bar{d}} + [G^{-1}J]_{\bar{a}} [G^{-1}J]_{\bar{b}} [G^{-1}J]_{\bar{c}} [G^{-1}J]_{\bar{d}} + (G^{-1})_{\bar{a}\bar{b}} [G^{-1}J]_{\bar{c}} [G^{-1}J]_{\bar{d}} \right).$$

##### 3.1.2 Replica representation of an inverse matrix

To obtain gaussian integrals we will use the "replica" representation the element  $(ij)$  of a matrix  $M$  of size  $D$ :

$$M_{ij}^{-1} = \lim_{n \rightarrow 0} \int \left( \prod_{\alpha=1}^n \prod_{i=1}^D d\eta_i^\alpha \right) \eta_i^1 \eta_j^1 \exp \left( -\frac{1}{2} \eta_i^\alpha M_{ij} \eta_j^\alpha \right). \quad (6)$$

Indeed, using the gaussian integral representation of the inverse of  $M$ ,

$$M_{ij}^{-1} = \mathcal{Z}^{-1} \int \left( \prod_{i=1}^D d\eta_i \right) \eta_i \eta_j \exp \left( -\frac{1}{2} \eta_i M_{ij} \eta_j \right),$$

$$\mathcal{Z} = \sqrt{\frac{(2\pi)^D}{\det M}}$$

$$= \int \left( \prod_{i=1}^D d\eta_i \right) \exp \left( -\frac{1}{2} \eta_i M_{ij} \eta_j \right).$$

Using the replica identity, we rewrite this as

$$M_{ij}^{-1} = \lim_{n \rightarrow 0} \mathcal{Z}^{n-1} \int \left( \prod_{i=1}^D d\eta_i \right) \eta_i \eta_j \exp \left( -\frac{1}{2} \eta_i M_{ij} \eta_j \right).$$

Renaming the integration variable of the integral on the left as  $\eta^1$  and the  $n - 1$  others as  $\eta^\alpha, \alpha \in \{2, n\}$ , we obtain expression (6).

### 3.2 The Random Feature model

In what follows, we will explicitly leave the indices of all the quantities used. We use the notation, called Einstein summation convention in physics, in which all repeated indices are summed but the sum is not explicitly written. Indices  $i \in \{1 \dots D\}$  are used to refer to the input dimension,  $h \in \{1 \dots P\}$  to refer to the hidden layer dimension and  $\mu \in \{1 \dots N\}$  to refer to the number of data points.

#### 3.2.1 With a single learner

In the random features model, the predictor can be computed explicitly:

$$\hat{\mathbf{a}} = \frac{1}{\sqrt{D}} \mathbf{y}_\mu^\top \mathbf{Z}_{\mu h} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)_{hh'}^{-1} \quad (7)$$

$$f(\mathbf{x}) = \hat{\mathbf{a}}_h \sigma \left( \frac{\Theta_{h'i} \mathbf{x}_i}{\sqrt{D}} \right) \quad (8)$$

$$= \mathbf{y}_\mu^\top \mathbf{Z}_{\mu h} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)_{hh'}^{-1} \sigma \left( \frac{\Theta_{h'i} \mathbf{x}_i}{\sqrt{D}} \right) / \sqrt{D}, \quad (9)$$

where

$$\mathbf{y}_\mu = f_d(\mathbf{X}_\mu) + \boldsymbol{\epsilon}_\mu, \quad (10)$$

$$\mathbf{Z}_{\mu h} = \frac{1}{\sqrt{D}} \sigma \left( \frac{1}{\sqrt{D}} \Theta_{hi} \mathbf{X}_{\mu i} \right). \quad (11)$$

Hence the test error can be computed as:

$$\mathcal{R}_{\text{RF}} = \mathbb{E}_{\mathbf{x}} \left[ \left( f_d(\mathbf{x}) - \mathbf{y}_\mu^\top \mathbf{Z}_{\mu h} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)_{hh'}^{-1} \sigma \left( \frac{\Theta_{h'i} \mathbf{x}_i}{\sqrt{D}} \right) / \sqrt{D} \right)^2 \right] \quad (12)$$

$$= \mathbb{E}_{\mathbf{x}} [f_d(\mathbf{x})^2] - 2 \mathbf{y}_\mu^\top \mathbf{Z}_{\mu h} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)_{hh'}^{-1} \mathbf{V}_h / \sqrt{D} \\ + \mathbf{y}_\mu^\top \mathbf{Z}_{\mu h} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)_{hh_1}^{-1} \mathbf{U}_{h_1 h_2} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)_{h_2 h'}^{-1} \mathbf{Z}_{h' \mu'}^\top \mathbf{y}_{\mu'} / D, \quad (13)$$

where

$$\mathbf{V}_h = \mathbb{E}_{\mathbf{x}} \left[ f_d(\mathbf{x}) \sigma \left( \frac{\langle \Theta_{hi} \mathbf{x}_i \rangle}{\sqrt{D}} \right) \right], \quad (14)$$

$$\mathbf{U}_{hh'} = \mathbb{E}_{\mathbf{x}} \left[ \sigma \left( \frac{\langle \Theta_{hi} \mathbf{x}_i \rangle}{\sqrt{D}} \right) \sigma \left( \frac{\Theta_{h'i} \mathbf{x}_i}{\sqrt{D}} \right) \right]. \quad (15)$$

#### 3.2.2 Ensembling over $K$ learners

When ensembling over  $K$  learners with independently sampled random feature vectors, the predictor becomes:

$$f(\mathbf{x}) = \frac{1}{K\sqrt{D}} \sum_k \mathbf{y}_\mu^\top \mathbf{Z}_{\mu h}^{(k)} \left( \mathbf{Z}^\top \mathbf{Z}^{(k)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)_{hh'}^{-1} \sigma \left( \frac{\Theta_{h'i}^{(k)} \mathbf{x}_i}{\sqrt{D}} \right), \quad (16)$$

where

$$\mathbf{Z}_{\mu h}^{(k)} = \frac{1}{\sqrt{D}} \sigma \left( \frac{1}{\sqrt{D}} \Theta_{hi}^{(k)} \mathbf{X}_{\mu i} \right). \quad (17)$$

The generalisation error is then given by:

$$\mathcal{R}_{\text{RF}} = \mathbb{E}_{\mathbf{x}} \left[ \left( f_d(\mathbf{x}) - \frac{1}{K} \sum_k \mathbf{y}^\top \mathbf{Z}^{(k)} \left( \mathbf{Z}^{\text{T}(k)} \mathbf{Z}^{(k)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \sigma \left( \frac{\langle \boldsymbol{\Theta}_{h'i}^{(k)} \mathbf{x}_i \rangle}{\sqrt{D}} \right) / \sqrt{D} \right)^2 \right] \quad (18)$$

$$\begin{aligned} &= \mathbb{E}_{\mathbf{x}} [f_d(\mathbf{x})^2] - \frac{2}{K} \sum_k \mathbf{y}^\top \mathbf{Z}^{(k)} \left( \mathbf{Z}^{\text{T}(k)} \mathbf{Z}^{(k)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \mathbf{V}^{(k)} / \sqrt{D} \\ &\quad + \frac{1}{K^2} \sum_{\substack{k \\ l \neq k}} \mathbf{y}^\top \mathbf{Z}^{(k)} \left( \mathbf{Z}^{\text{T}(k)} \mathbf{Z}^{(k)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \mathbf{U}^{(kl)} \left( \mathbf{Z}^{\text{T}(l)} \mathbf{Z}^{(l)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \mathbf{Z}^{(l)} \mathbf{T} \mathbf{y} / D, \end{aligned} \quad (19)$$

where

$$\mathbf{V}_h^{(k)} = \mathbb{E}_{\mathbf{x}} \left[ f_d(\mathbf{x}) \sigma \left( \frac{\langle \boldsymbol{\Theta}_{hi}^{(k)} \mathbf{x}_i \rangle}{\sqrt{D}} \right) \right], \quad (20)$$

$$\mathbf{U}_{hh'}^{(kl)} = \mathbb{E}_{\mathbf{x}} \left[ \sigma \left( \frac{\langle \boldsymbol{\Theta}_{hi}^{(k)} \mathbf{x}_i \rangle}{\sqrt{D}} \right) \sigma \left( \frac{\langle \boldsymbol{\Theta}_{h'i'}^{(l)} \mathbf{x}_{i'} \rangle}{\sqrt{D}} \right) \right]. \quad (21)$$

### 3.2.3 Equivalent Gaussian Covariate Model

It was shown in [4] that the random features model is equivalent, in the high-dimensional limit of Assumption 2, to a Gaussian covariate model in which the activation function  $\sigma$  is replaced as:

$$\sigma \left( \frac{\langle \boldsymbol{\Theta}_{hi}^{(k)} \mathbf{X}_{\mu i} \rangle}{\sqrt{D}} \right) \rightarrow \mu_0 + \mu_1 \frac{\langle \boldsymbol{\Theta}_{hi}^{(k)} \mathbf{X}_{\mu i} \rangle}{\sqrt{D}} + \mu_* \mathbf{W}_{\mu h}^{(k)}, \quad (22)$$

with  $\mathbf{W}^{(k)} \in \mathbb{R}^{N \times P}$ ,  $W_{\mu h}^{(k)} \sim \mathcal{N}(0, 1)$  and  $\mu_0, \mu_1$  and  $\mu_*$  defined in (2). To simplify the calculations, we take  $\mu_0 = 0$ , which amounts to adding a constant term to the activation function  $\sigma$ .

This powerful mapping allows to express the quantities  $\mathbf{U}, \mathbf{V}$ . We will not repeat their calculations here: the only difference here is  $\mathbf{U}^{kl}$ , which carries extra indices  $k, l$  due to the different initialization of the random features  $\boldsymbol{\Theta}^{(k)}$ . In our case,

$$\mathbf{U}_{hh'}^{(kl)} = \frac{\mu_1^2}{D} \boldsymbol{\Theta}_{hi}^{(k)} \boldsymbol{\Theta}_{h'i}^{(l)} + \mu_*^2 \delta_{kl} \delta_{hh'}. \quad (23)$$

Hence we can rewrite the test error as

$$\mathbb{E}_{\{\boldsymbol{\Theta}^{(k)}\}, \mathbf{X}, \varepsilon} [\mathcal{R}_{\text{RF}}] = F^2 (1 - 2\Psi_1^v) + \frac{1}{K} (F^2 \Psi_2^v + \tau^2 \Psi_3^v) + \left(1 - \frac{1}{K}\right) (F^2 \Psi_2^e + \tau^2 \Psi_3^e), \quad (24)$$

where  $\Psi_1, \Psi_2^v, \Psi_2^e, \Psi_3^v, \Psi_3^e$  are given by:

$$\begin{aligned} \Psi_1 &= \frac{1}{D} \text{Tr} \left[ \left( \frac{\mu_1}{D} \mathbf{X} \boldsymbol{\Theta}^{(1)\top} \right)^\top \mathbf{Z}^{(1)} \left( \mathbf{Z}^{(1)\top} \mathbf{Z}^{(1)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \right], \\ \Psi_2^v &= \frac{1}{D} \text{Tr} \left[ \left( \mathbf{Z}^{(1)\top} \mathbf{Z}^{(1)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \left( \frac{\mu_1^2}{D} \boldsymbol{\Theta}^{(1)} \boldsymbol{\Theta}^{(1)\top} + \mu_*^2 \mathbf{I}_N \right) \left( \mathbf{Z}^{(1)\top} \mathbf{Z}^{(1)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \mathbf{Z}^{(1)\top} \left( \frac{1}{D} \mathbf{X} \mathbf{X}^\top \right) \mathbf{Z}^{(1)} \right], \\ \Psi_3^v &= \frac{1}{D} \text{Tr} \left[ \left( \mathbf{Z}^{(1)\top} \mathbf{Z}^{(1)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \left( \frac{\mu_1^2}{D} \boldsymbol{\Theta}^{(1)} \boldsymbol{\Theta}^{(1)\top} + \mu_*^2 \mathbf{I}_N \right) \left( \mathbf{Z}^{(1)\top} \mathbf{Z}^{(1)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \mathbf{Z}^{(1)\top} \mathbf{Z}^{(1)} \right], \\ \Psi_2^e &= \frac{1}{D} \text{Tr} \left[ \left( \mathbf{Z}^{(1)\top} \mathbf{Z}^{(1)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \left( \frac{\mu_1^2}{D} \boldsymbol{\Theta}^{(1)} \boldsymbol{\Theta}^{(2)\top} \right) \left( \mathbf{Z}^{(2)\top} \mathbf{Z}^{(2)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \mathbf{Z}^{(2)\top} \left( \frac{1}{D} \mathbf{X} \mathbf{X}^\top \right) \mathbf{Z}^{(1)} \right], \\ \Psi_3^e &= \frac{1}{D} \text{Tr} \left[ \left( \mathbf{Z}^{(1)\top} \mathbf{Z}^{(1)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \left( \frac{\mu_1^2}{D} \boldsymbol{\Theta}^{(1)} \boldsymbol{\Theta}^{(2)\top} \right) \left( \mathbf{Z}^{(2)\top} \mathbf{Z}^{(2)} + \psi_1 \psi_2 \lambda \mathbf{I}_N \right)^{-1} \mathbf{Z}^{(2)\top} \mathbf{Z}^{(1)} \right]. \end{aligned}$$

### 3.3 Computation of the *vanilla* terms

To start with, let us compute the vanilla terms (those who carry a superscript  $v$ ), which involve a single instance of the random feature vectors. Note that these were calculated in [4] by evaluating the Stieljes transform of the



random matrices of which we need to calculate the trace. The replica method used here makes the calculation of the vanilla terms carry over easily to the the ensembling terms (superscript  $e$ ) and the divide and conquer term (superscript  $d$ ). To illustrate the calculation steps, we will calculate  $\Psi_3^v$ , then provide the results for  $\Psi_2^v$  and  $\Psi_1$ . In the vanilla terms, the two inverse matrices that appear are the same. Hence we use twice the replica identity (6), introducing  $2n$  replicas which all play the same role:

$$M_{ij}^{-1} M_{kl}^{-1} = \lim_{n \rightarrow 0} \int \left( \prod_{\alpha=1}^{2n} d\eta \right) \eta_i^1 \eta_j^1 \eta_k^2 \eta_l^2 \exp \left( -\frac{1}{2} \eta^\alpha M_{ij} \eta^\alpha \right). \quad (25)$$

The first step is to perform the averages, i.e. the Gaussian integrals, over the dataset  $X$ , the deterministic noise  $W$  induced by the non-linearity of the activation function and the random features  $\Theta$ .

### 3.3.1 Averaging over the dataset

Replacing the activation function by its Gaussian covariate equivalent model and using (25), the term  $\Psi_3$  can be expanded as:

$$\begin{aligned} \Psi_3^v &= \frac{1}{D} \left[ \mathbf{Z}_{\mu h} (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)_{hh_1}^{-1} \left( \frac{\mu_1^2}{D} \Theta_{h_1 i} \Theta_{h_2 i} + \mu_\star^2 \delta_{h_1 h_2} \right) (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)_{h_2 h'}^{-1} \mathbf{Z}_{h' \mu} \right] \\ &= \frac{1}{D} \left( \frac{\mu_1^2}{D} \Theta_{h_1 i} \Theta_{h_2 i} + \mu_\star^2 \delta_{h_1 h_2} \right) [\mathbf{Z}_{\mu h} \mathbf{Z}_{h' \mu}] \int \left( \prod_{\alpha=1}^{2n} d\eta^\alpha \right) \eta_h^1 \eta_{h_1}^1 \eta_{h_2}^2 \eta_{h'}^2 \exp \left( -\frac{1}{2} \eta_h^\alpha (\mathbf{Z}^\top \mathbf{Z} + \psi_1 \psi_2 \lambda \mathbf{I}_N)_{hh'} \eta_{h'}^\alpha \right) \\ &= \frac{1}{D^2} \left( \frac{\mu_1^2}{D} \Theta_{h_1 i} \Theta_{h_2 i} + \mu_\star^2 \delta_{h_1 h_2} \right) \left( \frac{\mu_1}{\sqrt{D}} \Theta \mathbf{X} + \mu_\star \mathbf{W} \right)_{h\mu} \left( \frac{\mu_1}{\sqrt{D}} \Theta \mathbf{X} + \mu_\star \mathbf{W} \right)_{h'\mu} \\ &\int \left( \prod_{\alpha=1}^{2n} d\eta^\alpha \right) \eta_h^1 \eta_{h_1}^1 \eta_{h_2}^2 \eta_{h'}^2 \exp \left( -\frac{1}{2} \eta_h^\alpha \left( \frac{1}{D} \left( \frac{\mu_1}{\sqrt{D}} \Theta \mathbf{X} + \mu_\star \mathbf{W} \right)_{h\mu} \left( \frac{\mu_1}{\sqrt{D}} \Theta \mathbf{X} + \mu_\star \mathbf{W} \right)_{h'\mu} + \psi_1 \psi_2 \lambda \delta_{hh'} \right) \eta_{h'}^\alpha \right). \end{aligned}$$

Now, we introduce  $\lambda_i^\alpha := \frac{1}{\sqrt{P}} \eta_h^\alpha \Theta_{hi}$ , and enforce this relation using the Fourier representation of the delta-function:

$$1 = \int d\lambda_i^\alpha d\hat{\lambda}_i^\alpha e^{i\hat{\lambda}_i^\alpha (\sqrt{P}\lambda_i^\alpha - \eta_h^\alpha \Theta_{hi})}. \quad (26)$$

The average over the dataset  $\mathbf{X}_{\mu i}$  has the form of (5) with:

$$(G_X)_{\mu\mu', ii'} = \delta_{\mu\mu'} \left( \delta_{ii'} + \frac{\mu_1^2 \psi_1}{D} \lambda_i^\alpha \lambda_{i'}^\alpha \right), \quad (27)$$

$$(J_X)_{\mu, i} = \frac{\mu_1 \mu_\star \sqrt{\psi_1}}{D} \sum_{\alpha\alpha} \lambda_i^\alpha \eta_h^\alpha \mathbf{W}_{\mu h}. \quad (28)$$

Using formulae (5), we obtain:

$$\begin{aligned} \Psi_3^v &= \frac{N}{D^2} \int \left( \prod_{\alpha=1}^{2n} d\eta^\alpha d\lambda^\alpha d\hat{\lambda}^\alpha \right) \left( \frac{\mu_1^2 P}{D} \lambda_i^1 \lambda_i^2 + \mu_\star^2 \eta_i^1 \eta_i^2 \right) \\ &\left[ \mu_\star^2 \eta_h^1 \mathbf{W}_h \mathbf{W}_{h'} \eta_{h'}^2 + \mu_1^2 \psi_1 \lambda_i^1 \lambda_{i'}^2 \left( (G_X^{-1})_{ii'} + (G_X^{-1} J_X)_i (G_X^{-1} J_X)_{i'} \right) + 2\mu_1 \mu_\star \sqrt{\psi_1} \lambda_i^1 \mathbf{W}_h \eta_h^2 (G_X^{-1} J_X)_i \right] \\ &\exp \left( -\frac{n}{2} \log \det(G_X) - \frac{1}{2} \eta_h^\alpha \left( \frac{\mu_\star^2}{D} \mathbf{W}_h \mathbf{W}_{h'} \delta_{\alpha\beta} - \frac{\mu_1^2 \mu_\star^2 \psi_1}{D^2} \mathbf{W}_h \lambda_i^\alpha (G_X)_{ii'}^{-1} \lambda_{i'}^\beta \mathbf{W}_{h'} + \psi_1 \psi_2 \lambda \delta_{hh'} \right) \eta_{h'}^\beta + i\hat{\lambda}_i^\alpha (\sqrt{P}\lambda_i^\alpha - \eta_h^\alpha \Theta_{hi}) \right) \end{aligned}$$

Note that due to with a slight abuse of notation we got rid of indices  $\mu$ , which all sum up trivially to give a global factor  $N$ .

### 3.3.2 Averaging over the deterministic noise

The expectation over the deterministic noise  $\mathbf{W}_h$  is a Gaussian integral of the form (5) with:

$$[G_W]_{hh'} = \delta_{hh'} + \frac{\mu_\star^2}{D} \eta_h^\alpha A^{\alpha\beta} \eta_{h'}^\beta, \quad (29)$$

$$[J_W]_h = 0, \quad (30)$$

$$A^{\alpha\beta} = \delta_{\alpha\beta} - \mu_1^2 \psi_1 \frac{1}{D} \sum_{i,j} \lambda_i^\alpha [G_X^{-1}]_{ij} \lambda_j^\beta. \quad (31)$$

Note that the prefactor involves, constant, linear and quadratic terms in  $\mathbf{W}$  since:

$$(G_X^{-1} J_X)_i = \frac{\mu_1 \mu_\star \sqrt{\psi_1}}{D} [\eta^\alpha \mathbf{W}] [G_X^{-1} \lambda^\alpha]_i.$$

Thus, one obtains:

$$\begin{aligned} \Psi_3^v = & \frac{\psi_2}{D} \int \left( \prod_{\alpha=1}^{2n} d\eta^\alpha d\lambda^\alpha d\hat{\lambda}^\alpha \right) (\mu_1^2 \psi_1 \lambda_i^1 \lambda_i^2 + \mu_\star^2 \eta_i^1 \eta_i^2) [\mu_\star^2 [\eta^1 (G_W^{-1}) \eta^2] + \mu_1^2 \psi_1 [\lambda^1 H_W \lambda^2] + 2\mu_1^2 \mu_\star^2 \psi_1 [\lambda^1 S_W \eta^2]] \\ & \exp \left( -\frac{n}{2} \log \det(G_X) - \frac{n}{2} \log \det(G_W) - \frac{1}{2} \psi_1 \psi_2 \lambda \sum (\eta_h^\alpha)^2 + i \hat{\lambda}_i^\alpha (\sqrt{P} \lambda_i^\alpha - \eta_h^\alpha \Theta_{hi}) \right), \end{aligned}$$

with

$$(H_W)_{ij} = (G_X^{-1})_{ij} + \frac{\mu_1^2 \mu_\star^2 \psi_1}{D^2} [\eta^\alpha (G_W^{-1}) \eta^\beta] [G_X^{-1} \lambda^\alpha]_i [G_X^{-1} \lambda^\beta]_j, \quad (32)$$

$$(S_W)_{ih} = \frac{1}{D} [G_X^{-1} \lambda^\alpha]_i [G_W^{-1} \eta^\alpha]_h. \quad (33)$$

### 3.3.3 Averaging over the random feature vectors

The expectation over the random feature vectors  $\Theta_{hi}$  is a Gaussian integral of the form (5) with:

$$[G_\Theta]_{hh',ii'} = \delta_{hh',ii'}, \quad (34)$$

$$[J_\Theta]_{hi} = -i \hat{\lambda}_i^\alpha \eta_h^\alpha. \quad (35)$$

Performing this integration results in:

$$\begin{aligned} \Psi_3^v = & \frac{\psi_2}{D} \int \left( \prod_{\alpha=1}^{2n} d\eta^\alpha d\lambda^\alpha d\hat{\lambda}^\alpha \right) (\mu_1^2 \psi_1 \lambda_i^1 \lambda_i^2 + \mu_\star^2 \eta_i^1 \eta_i^2) [\mu_\star^2 [\eta^1 (G_W^{-1}) \eta^2] + \mu_1^2 \psi_1 [\lambda^1 H_W \lambda^2] + 2\mu_1^2 \mu_\star^2 \psi_1 [\lambda^1 S_W \eta^2]] \\ & \exp \left( -\frac{n}{2} \log \det(G_X) - \frac{n}{2} \log \det(G_W) - \frac{1}{2} \psi_1 \psi_2 \lambda \sum (\eta_h^\alpha)^2 - \frac{1}{2} \eta_h^\alpha \eta_h^\beta \hat{\lambda}_i^\alpha \hat{\lambda}_i^\beta + i \sqrt{P} \hat{\lambda}_i^\alpha \lambda_i^\alpha \right). \end{aligned}$$

### 3.3.4 Expression of the action and the prefactor

To complete the computation we integrate with respect to  $\hat{\lambda}_i^\alpha$ , using again formulae (5):

$$[G_\lambda]_{ii'}^{\alpha\beta} = \delta^{ii'} \eta_h^\alpha \eta_h^\beta, \quad (36)$$

$$[J_\lambda]_i^\alpha = i \sqrt{P} \lambda_i^\alpha. \quad (37)$$

This yields the final expression of the term:

$$\begin{aligned} \Psi_3^v = & \frac{\psi_2}{D} \int \left( \prod_{\alpha=1}^{2n} d\eta^\alpha \right) \left( \prod_{\alpha=1}^{2n} d\lambda^\alpha \right) (\mu_1^2 \psi_1 \lambda_i^1 \lambda_i^2 + \mu_\star^2 \eta_i^1 \eta_i^2) [\mu_\star^2 [\eta^1 (G_W^{-1}) \eta^2] + \mu_1^2 \psi_1 [\lambda^1 H_W \lambda^2] + 2\mu_1^2 \mu_\star^2 \psi_1 [\lambda^1 S_W \eta^2]] \\ & \exp \left( -\frac{n}{2} \log \det(G_X) - \frac{n}{2} \log \det(G_W) - \frac{D}{2} \log \det(G_\lambda) - \frac{1}{2} \psi_1 \psi_2 \lambda \sum (\eta_h^\alpha)^2 - \frac{P}{2} \lambda_i^\alpha (G_\lambda^{-1})_{ii'}^{\alpha\beta} \lambda_i^\beta \right). \end{aligned}$$

The above may be written as

$$\Psi_3^v = \int \left( \prod d\eta \right) \left( \prod d\lambda \right) P_{\Psi_3^v} [\eta, \lambda] \exp \left( -\frac{D}{2} S^v [\eta, \lambda] \right), \quad (38)$$

with the *prefactor*  $P_{\Psi_3^v}$  and the *action*  $S^v$  defined as:

$$P_{\Psi_3^v} [\eta, \lambda] := \frac{\psi_2}{D} \left( \mu_1^2 \psi_1 \lambda_i^1 \lambda_i^2 + \mu_\star^2 \eta_i^1 \eta_i^2 \right) \left[ \mu_\star^2 [\eta^1 (G_W^{-1}) \eta^2] + \mu_1^2 \psi_1 [\lambda^1 H_W \lambda^2] + 2\mu_1^2 \mu_\star^2 \psi_1 [\lambda^1 S_W \eta^2] \right],$$

$$S^v [\eta, \lambda] := \psi_2 \log \det(G_X) + \psi_2 \log \det(G_W) + \log \det(G_\lambda) + \frac{1}{D} \psi_1 \psi_2 \lambda \sum (\eta_h^\alpha)^2 + \frac{P}{D} \left( \lambda_i^\alpha (G_\lambda^{-1})_{ii'}^{\alpha\beta} \lambda_{i'}^\beta \right).$$

### 3.3.5 Expression of the action and the prefactor in terms of order parameters

Here we see that we have a factor  $D \rightarrow \infty$  in the exponential part, which can be estimated using the saddle point method. Before doing so, we introduce the following order parameters using the Fourier representation of the delta-function:

$$1 = \int dQ_{\alpha\beta} d\hat{Q}_{\alpha\beta} e^{\hat{Q}_{\alpha\beta} (PQ_{\alpha\beta} - \eta_h^\alpha \eta_h^\beta)}, \quad (39)$$

$$1 = \int dR_{\alpha\beta} d\hat{R}_{\alpha\beta} e^{\hat{R}_{\alpha\beta} (DR_{\alpha\beta} - \lambda_i^\alpha \lambda_i^\beta)}. \quad (40)$$

This allows to rewrite the prefactor only in terms of  $Q, R$ : for example,

$$\mu_1^2 \psi_1 \lambda_i^1 \lambda_i^2 + \mu_\star^2 \eta_i^1 \eta_i^2 = \psi_1 D (\mu_1^2 R^{12} + \mu_\star^2 Q^{12}).$$

To do this, there are two key quantities we need to calculate:  $\lambda G_X^{-1} \lambda$  and  $\eta G_W^{-1} \eta$ . To calculate both, we note that  $G_X$  and  $G_W$  are both of the form  $\mathbf{I} + \mathbf{X}$ , therefore their inverse may be calculated using their series representation. The result is:

$$[M_X]^{\alpha\beta} := \frac{1}{D} \lambda^\alpha G_X^{-1} \lambda^\beta = [R(I + \mu_1^2 \psi_1 R)^{-1}]^{\alpha\beta}, \quad (41)$$

$$[M_W]^{\alpha\beta} := \frac{1}{P} \eta^\alpha (G_W^{-1}) \eta^\beta = [Q(I + \mu_\star^2 \psi_1 A Q)^{-1}]^{\alpha\beta}. \quad (42)$$

Using the above, we deduce:

$$\lambda^1 H_W \lambda^2 = D M_X^{12} + P \mu_1^2 \mu_\star^2 \psi_1 [M_X M_W M_X]^{12}, \quad (43)$$

$$\lambda^1 S_W \eta^2 = P [M_X M_W]^{12}. \quad (44)$$

The integrals over  $\eta, \lambda$  become simple Gaussian integrals with covariance matrices given by  $\hat{Q}, \hat{R}$ , yielding:

$$1 = \int dQ_{\alpha\beta} d\hat{Q}_{\alpha\beta} e^{-\frac{\psi_1 D}{2} (\log \det \hat{Q} - 2 \text{Tr} Q \hat{Q})}, \quad (45)$$

$$1 = \int dR_{\alpha\beta} d\hat{R}_{\alpha\beta} e^{-\frac{D}{2} (\log \det \hat{R} - 2 \text{Tr} R \hat{R})}. \quad (46)$$

The next step is to take the saddle point with respect to the auxiliary variables  $\hat{Q}$  and  $\hat{R}$  in order to eliminate them:

$$\frac{\partial S^v}{\partial \hat{Q}_{\alpha\beta}} = \psi_1 \left( \hat{Q}^{-1} - 2Q \right) = 0 \Rightarrow \hat{Q} = \frac{1}{2} Q^{-1}, \quad (47)$$

$$\frac{\partial S^v}{\partial \hat{R}_{\alpha\beta}} = \left( \hat{R}^{-1} - 2R \right) = 0 \Rightarrow \hat{R} = \frac{1}{2} R^{-1}. \quad (48)$$

One finally obtains that:

$$\Psi_3^v = \int \left( \prod dQ \right) \left( \prod dR \right) P_{\Psi_3^v} [Q, R] \exp \left( -\frac{D}{2} S^v [Q, R] \right), \quad (49)$$

With:

$$P_{\Psi_3^v} [Q, R] = D \psi_1^2 \psi_2 (\mu_1^2 R^{12} + \mu_\star^2 Q^{12}) \left[ \mu_\star^2 [M_W^v]^{12} + \mu_1^2 \left( M_X^{12} + \mu_1^2 \mu_\star^2 \psi_1^2 [M_X M_W^v M_X]^{12} \right) + 2\mu_1^2 \mu_\star^2 \psi_1 [M_X M_W^v]^{12} \right],$$

$$S^v [Q, R] = \psi_2 \log \det(G_X) + \psi_2 \log \det(G_W) + \psi_1^2 \psi_2 \lambda \text{Tr} Q + \text{Tr} (R Q^{-1}) + (1 - \psi_1) \log \det Q - \log \det R. \quad (50)$$

### 3.3.6 Saddle point equations

The aim is now to use the saddle point method in order to evaluate the integrals over the order parameters. Thus, one looks for  $R$  and  $Q$  solutions to the equations:

$$\frac{\partial S^v}{\partial Q_{\alpha\beta}} = 0, \quad \frac{\partial S^v}{\partial R_{\alpha\beta}} = 0 \quad \forall \alpha, \beta = 1, \dots, 2n.$$

To solve the above, it is common to make a *replica symmetric ansatz*. In this case, we assume that the solutions to the saddle points equations take the form:

$$Q = \begin{bmatrix} q & \tilde{q} & \cdots & \tilde{q} \\ \tilde{q} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \tilde{q} \\ \tilde{q} & \cdots & \tilde{q} & q \end{bmatrix}, \quad R = \begin{bmatrix} r & \tilde{r} & \cdots & \tilde{r} \\ \tilde{r} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \tilde{r} \\ \tilde{r} & \cdots & \tilde{r} & r \end{bmatrix} \quad (51)$$

The action takes the following form:

$$\begin{aligned} S^v(q, r, \tilde{q}, \tilde{r}) &= 2n(S_0(q, r) + S_1^v(q, r, \tilde{q}, \tilde{r})) \\ S_0(q, r) &= \lambda\psi_1^2\psi_2q + \psi_2 \log\left(\frac{\mu_*^2\psi_1q}{\mu_1^2\psi_1r+1} + 1\right) + \frac{r}{q} + (1-\psi_1)\log(q) + \psi_2 \log(\mu_1^2\psi_1r+1) - \log(r) \\ S_1^v(q, r, \tilde{q}, \tilde{r}) &= f^v(q, r)\tilde{q} + g^v(q, r)\tilde{r} \\ f^v(q, r) &= \lambda\psi_1^2\psi_2 + \frac{\mu_*^2\psi_1\psi_2}{\mu_*^2\psi_1q + \mu_1^2\psi_1r + 1} + \frac{1-\psi_1}{q} - \frac{r}{q^2} \\ g^v(q, r) &= -\frac{\mu_*^2\mu_1^2\psi_1^2\psi_2q}{(\mu_1^2\psi_1r+1)(\mu_*^2\psi_1q + \mu_1^2\psi_1r + 1)} + \frac{\mu_1^2\psi_1\psi_2}{\mu_1^2\psi_1r+1} + \frac{1}{q} - \frac{1}{r}. \end{aligned} \quad (52)$$

### 3.3.7 Fluctuations around the saddle point

We introduce the following notations:

$$\begin{aligned} [\nabla_{\mathbf{T}} F(\mathbf{T}_*)]_{\alpha\beta} &= \frac{\partial F}{\partial \mathbf{T}_{\alpha\beta}} \Big|_{\mathbf{T}_*}, \\ [H_{\mathbf{T}} F(\mathbf{T}_*)]_{\alpha\beta, \gamma\delta} &= \left[ \begin{array}{cc} \frac{\partial^2 F}{\partial Q_{\alpha\beta} \partial Q_{\gamma\delta}} & \frac{\partial^2 F}{\partial Q_{\alpha\beta} \partial R_{\gamma\delta}} \\ \frac{\partial^2 F}{\partial R_{\alpha\beta} \partial Q_{\gamma\delta}} & \frac{\partial^2 F}{\partial R_{\alpha\beta} \partial R_{\gamma\delta}} \end{array} \right] \Big|_{\mathbf{T}_*}, \\ H[F] &= \left[ \begin{array}{cccc} \frac{\partial F}{\partial q} & \frac{\partial F}{\partial r} & \frac{\partial F}{\partial \tilde{q}} & \frac{\partial F}{\partial \tilde{r}} \\ \frac{\partial q}{\partial q} & \frac{\partial q}{\partial r} & \frac{\partial q}{\partial \tilde{q}} & \frac{\partial q}{\partial \tilde{r}} \\ \frac{\partial r}{\partial q} & \frac{\partial r}{\partial r} & \frac{\partial r}{\partial \tilde{q}} & \frac{\partial r}{\partial \tilde{r}} \\ \frac{\partial \tilde{q}}{\partial q} & \frac{\partial \tilde{q}}{\partial r} & \frac{\partial \tilde{q}}{\partial \tilde{q}} & \frac{\partial \tilde{q}}{\partial \tilde{r}} \\ \frac{\partial \tilde{r}}{\partial q} & \frac{\partial \tilde{r}}{\partial r} & \frac{\partial \tilde{r}}{\partial \tilde{q}} & \frac{\partial \tilde{r}}{\partial \tilde{r}} \end{array} \right] \Big|_{\substack{q=q^* \\ r=r^* \\ \tilde{q}=0 \\ \tilde{r}=0}}. \end{aligned}$$

**Proposition** Let  $q^*$  and  $r^*$  be the solutions of the fixed point equation for the function  $S_0 : (q, r) \mapsto \mathbb{R}$  defined in (52):

$$\begin{cases} \frac{\partial S_0(q, r)}{\partial q} = 0 \\ \frac{\partial S_0(q, r)}{\partial r} = 0. \end{cases}$$

Then we have that

$$\Psi_3^v = \frac{1}{D} \text{Tr} \left[ H[S^v]^{-1} H[P_{\Psi_3^v}] \right]. \quad (53)$$

**Sketch of proof** Solving the saddle point equations:

$$\begin{cases} \frac{\partial S^v(q,r,\tilde{q},\tilde{r})}{\partial q} = 0 \\ \frac{\partial S^v(q,r,\tilde{q},\tilde{r})}{\partial r} = 0 \\ \frac{\partial S^v(q,r,\tilde{q},\tilde{r})}{\partial \tilde{q}} = 0 \\ \frac{\partial S^v(q,r,\tilde{q},\tilde{r})}{\partial \tilde{r}} = 0 \end{cases},$$

one finds  $\tilde{q} = \tilde{r} = 0$ , which is problematic because the prefactor vanishes:  $P_{\Psi_3} \propto \mu_1^2 \tilde{q} + \mu_2^2 \tilde{r}$ .

Therefore we must go beyond the saddle point contribution to obtain a non zero result, i.e. we have to examine the quadratic fluctuations around the saddle point. To do so we perform a second-order expansion of the action (50) as a function of  $Q$  and  $R$ :

$$\begin{aligned} P_{\Psi_3^v}(\mathbf{T}) &\approx P_{\Psi_3^v}(\mathbf{T}_\star) + (\mathbf{T} - \mathbf{T}_\star)^\top \nabla P_{\Psi_3^v}(\mathbf{T}_\star) + \frac{1}{2} (\mathbf{T} - \mathbf{T}_\star)^\top H_{\mathbf{T}} [P_{\Psi_3^v}(\mathbf{T}_\star)] (\mathbf{T} - \mathbf{T}_\star), \\ S^v(\mathbf{T}) &\approx S_0(\mathbf{T}_\star) + \frac{1}{2} (\mathbf{T} - \mathbf{T}_\star)^\top H_{\mathbf{T}} [S^v(\mathbf{T}_\star)] (\mathbf{T} - \mathbf{T}_\star). \end{aligned}$$

Computing the second derivative of (50), it is easy to show that:

$$\begin{aligned} [H_{\mathbf{T}} [S^v(\mathbf{T})]]_{\alpha\beta,\gamma\delta} &= [H_{\mathbf{T}} [S^v(\mathbf{T})]]_{\alpha\beta} (\delta_{\alpha\gamma} \delta_{\beta\delta} + \delta_{\alpha\delta} \delta_{\beta\gamma}), \\ [H_{\mathbf{T}} [P_{\Psi_3^v}(\mathbf{T})]]_{\alpha\beta,\gamma\delta} &= [H_{\mathbf{T}} [P_{\Psi_3^v}(\mathbf{T})]]_{\alpha\beta} (\delta_{\alpha\gamma} \delta_{\beta\delta} + \delta_{\alpha\delta} \delta_{\beta\gamma}), \end{aligned}$$

where

$$\begin{aligned} H_{\mathbf{T}} [F]_{\alpha\beta} &= \begin{bmatrix} \frac{\partial^2 F}{\partial Q_{\alpha\beta} \partial Q_{\alpha\beta}} & \frac{\partial^2 F}{\partial Q_{\alpha\beta} \partial R_{\alpha\beta}} \\ \frac{\partial^2 F}{\partial R_{\alpha\beta} \partial Q_{\alpha\beta}} & \frac{\partial^2 F}{\partial R_{\alpha\beta} \partial R_{\alpha\beta}} \end{bmatrix} \\ &= \frac{1}{2n} \delta_{\alpha\beta} \begin{bmatrix} \frac{\partial^2 F}{\partial q \partial q} & \frac{\partial^2 F}{\partial q \partial r} \\ \frac{\partial^2 F}{\partial r \partial q} & \frac{\partial^2 F}{\partial r \partial r} \end{bmatrix} + \frac{2}{2n(2n-1)} (1 - \delta_{\alpha\beta}) \begin{bmatrix} \frac{\partial^2 F}{\partial \tilde{q} \partial \tilde{q}} & \frac{\partial^2 F}{\partial \tilde{q} \partial \tilde{r}} \\ \frac{\partial^2 F}{\partial \tilde{r} \partial \tilde{q}} & \frac{\partial^2 F}{\partial \tilde{r} \partial \tilde{r}} \end{bmatrix}. \end{aligned}$$

Hence,

$$\begin{aligned} \Psi_3^v &= \lim_{n \rightarrow 0} \int d\mathbf{T} P_{\Psi_3^v}(\mathbf{T}) \exp^{-\frac{D}{2} S^v(\mathbf{T})}, \\ &= \lim_{n \rightarrow 0} \underbrace{P_{\Psi_3^v}(\mathbf{T}_\star)}_0 + \nabla P_{\Psi_3^v}(\mathbf{T}_\star)_{\alpha\beta} \underbrace{\int d\mathbf{T} (\mathbf{T} - \mathbf{T}_\star)_{\alpha\beta} \exp\left(-\frac{D}{2} S^v(\mathbf{T})\right)}_0 \\ &\quad + \frac{1}{2} H_{\mathbf{T}} [P_{\Psi_3^v}(\mathbf{T}_\star)]_{\alpha\beta} \int d\mathbf{T} (\mathbf{T} - \mathbf{T}_\star)_{\alpha\beta}^2 \exp\left(-\frac{D}{2} S^v(\mathbf{T})\right) \\ &= \lim_{n \rightarrow 0} \frac{1}{2} H_{\mathbf{T}} [P_{\Psi_3^v}(\mathbf{T}_\star)]_{\alpha\beta} e^{-\frac{D}{2} S^v(\mathbf{T}_\star)} \int d\mathbf{T} (\mathbf{T} - \mathbf{T}_\star)_{\alpha\beta}^2 e^{-\frac{D}{2} \sum_{\alpha\beta} (\mathbf{T} - \mathbf{T}_\star)_{\alpha\beta}^2} H_{\mathbf{T}} [S^v(\mathbf{T}_\star)]_{\alpha\beta} \\ &= \lim_{n \rightarrow 0} \frac{1}{2} \frac{(2\pi)^n}{\det H_{\mathbf{T}} [S^v(\mathbf{T}_\star)]} e^{-\frac{D}{2} S^v(\mathbf{T}_\star)} \frac{2}{D} H_{\mathbf{T}} [P_{\Psi_3^v}(\mathbf{T}_\star)]_{\alpha\beta} H_{\mathbf{T}} [S^v(\mathbf{T}_\star)]_{\alpha\beta}^{-1} \\ &= \frac{1}{D} \text{Tr} \left[ H [S^v]^{-1} H [P_{\Psi_3^v}] \right] \end{aligned}$$

In the last step, we used the fact that:

$$\begin{aligned} \lim_{n \rightarrow 0} e^{-\frac{D}{2} S^v(\mathbf{T}_\star)} &= \lim_{n \rightarrow 0} e^{-n D S_0(q_\star, r_\star)} = 1, \\ \lim_{n \rightarrow 0} \det H_{\mathbf{T}} [S^v(\mathbf{T}_\star)] &= 1. \end{aligned}$$

The last equality follows from the fact that for a matrix of size  $n \times n$  of the form  $M_{\alpha\beta} = a\delta_{\alpha\beta} + b(1 - \delta_{\alpha\beta})$ , we have

$$\det M = (a - b)^n \left( 1 + \frac{nb}{a - b} \right) \xrightarrow{n \rightarrow 0} 1.$$

### 3.3.8 Expression of the vanilla terms

Using the above procedure, one can compute the terms  $\Psi_1, \Psi_2^g, \Psi_3^g$  of (24): for each of these terms, the action is the same as in (50), and the prefactors can be obtained as:

$$\begin{aligned}
P_{\Psi_1} [Q, R] &= \mu_1^2 \psi_1 \psi_2 \left( M_X^{11} + \mu_1^2 \mu_\star^2 \psi_1^2 (M_X M_W M_X)^{11} + \mu_\star^2 \mu_1 \psi_1 (M_X M_W)^{11} \right), \\
P_{\Psi_2^g} [Q, R] &= D \psi_1^2 \psi_2 (\mu_1^2 R^{12} + \mu_\star^2 Q^{12}) (\mu_1^2 \psi_2 P_{XX} - 2\mu_1^2 \mu_\star^2 \psi_1 \psi_2 P_{XW} + \mu_\star^2 P_{WW}), \\
P_{\Psi_3^g} [Q, R] &= D \psi_1^2 \psi_2 (\mu_1^2 R^{12} + \mu_\star^2 Q^{12}) \left[ \mu_\star^2 [M_W^v]^{12} + \mu_1^2 \left( M_X^{12} + \mu_1^2 \mu_\star^2 \psi_1^2 [M_X M_W^v M_X]^{12} \right) + 2\mu_1^2 \mu_\star^2 \psi_1 [M_X M_W^v]^{12} \right], \\
P_{XX} &= N_X^{12} + \frac{1}{\psi_2} M_X^{12} + 2(\mu_1 \mu_\star \psi_1)^2 [N_X M_W M_X]^{12} + \frac{(\mu_1 \mu_\star \psi_1)^2}{\psi_2} [M_X M_W M_X]^{12} + (\mu_1 \mu_\star \psi_1)^4 [M_X M_W N_X M_W M_X]^{12}, \\
P_{XW} &= [N_X M_W]^{12} + \frac{1}{\psi_2} [M_X M_W]^{12} + (\mu_1 \mu_\star \psi_1)^2 [M_X M_W N_X M_W]^{12}, \\
P_{WW} &= M_W^{12} + \psi_2 (\mu_1 \mu_\star \psi_1)^2 [M_W N_X M_W]^{12}.
\end{aligned}$$

Where a new term appears:

$$[N_X]^{\alpha\beta} = [R(I + \mu_1^2 \psi_1 R)^{-2}]^{\alpha\beta}. \quad (54)$$

## 3.4 Computation of the *ensembling* terms

### 3.4.1 Expression of the action and the prefactor

In the *ensembling* terms, the two inverse matrices are different, hence one has to introduce two distinct replica variables. We distinguish them by the use of an extra index  $a \in \{1, 2\}$ , denoted in brackets in order not to be confused with the replica indices  $\alpha$ .

$$\left[ M^{(1)} \right]_{ij}^{-1} \left[ M^{(2)} \right]_{kl}^{-1} = \lim_{n \rightarrow 0} \int \left( \prod_{\alpha=1}^n \prod_{a=1}^2 d\eta^{\alpha(a)} \right) \eta_i^{1(1)} \eta_j^{1(1)} \eta_k^{1(2)} \eta_l^{1(2)} \exp \left( -\frac{1}{2} \sum_{(a)} \eta^{\alpha(a)} M_{ij}^{(a)} \eta^{\alpha(a)} \right). \quad (55)$$

Calculations of the Gaussian integrals follow through in a very similar way as for the *vanilla* terms. The matrices appearing in the process are:

$$(G_X^e)_{ii'} = \delta_{ii'} + \frac{\mu_1^2 \psi_1}{D} \sum_{(a)\alpha} \lambda_i^{\alpha(a)} \lambda_{i'}^{\alpha(a)}, \quad (56)$$

$$(J_X^e)_i = \frac{\mu_1 \mu_\star \sqrt{\psi_1}}{D^2} \sum_{\alpha a} \lambda_i^{\alpha(a)} \eta_h^{\alpha(a)} W_h^{(a)}, \quad (57)$$

$$(G_W^e)_{hh'}^{(ab)} = \delta_{hh'}^{ab} + \frac{\mu_\star^2}{D} \sum_{\alpha\beta} \eta_h^{\alpha(a)} A_e^{\alpha\beta, ab} \eta_{h'}^{\beta(b)}, \quad (58)$$

$$(J_W^e)_h = 0 \quad (59)$$

$$(G_\Theta^e)_{hh', ii'}^{(ab)} = \delta_{hh', ii'}^{ab} \quad (60)$$

$$(J_\Theta^e)_{hi}^{(a)} = -i \sum_{\alpha} \hat{\lambda}_i^{\alpha(a)} \eta_h^{\alpha(a)}, \quad (61)$$

$$(G_\lambda^e)_{ii'}^{\alpha\beta, (ab)} = 2\delta^{ab} \delta^{ii'} \eta_h^{\alpha(a)} \eta_h^{\beta(b)}, \quad (62)$$

$$(J_\lambda^e)_i^{\alpha(a)} = i\sqrt{P} \lambda_i^{\alpha(a)}. \quad (63)$$

$$(64)$$

with

$$A_e^{\alpha\beta,(ab)} = \delta_{ab}\delta_{\alpha\beta} - \mu_1^2\psi_1 \frac{1}{D} \sum_{i,j} \lambda_i^{\alpha(a)} [G_X^{e-1}]_{ij} \lambda_j^{\beta(b)}, \quad (65)$$

$$[H_W^e]_{ii'} = [G_X^{e-1}]_{ii'} + \sum_{\alpha\beta,ab} \left[ G_X^{e-1} \lambda^{\alpha(a)} \right]_i \left[ G_X^{e-1} \lambda^{\beta(b)} \right]_{i'} \left[ \eta^{\alpha(a)} [G_W^{e-1}]^{(ab)} \eta^{\beta(b)} \right], \quad (66)$$

$$[S_W^e]_{ih} = \frac{1}{D} \sum_{\alpha,a} \left[ G_X^{e-1} \lambda^{\alpha(a)} \right]_i \left[ G_W^{e-1} \eta^{\alpha(a)} \right]_h. \quad (67)$$

Starting with the computation of  $\Psi_3^e$  in order to illustrate the method used, the prefactor  $P_{\Psi_3^e}$  and the action are  $S^e$  are given by:

$$P_{\Psi_3^e} [\eta, \lambda] := \frac{\mu_1^2\psi_1\psi_2}{D} \lambda_i^{1(1)} \lambda_i^{1(2)} \left[ \mu_*^2 \left[ \eta^{1(1)} (G_W^{e-1})^{(12)} \eta^{1(2)} \right] + \mu_1^2\psi_1 \left[ \lambda^{1(1)} H_W^e \lambda^{1(2)} \right] + 2\mu_1^2\mu_*^2\psi_1 \left[ \lambda^{1(1)} S_W^e \eta^{1(2)} \right] \right], \quad (68)$$

$$S^e [\eta, \lambda] := \psi_2 \log \det(G_X^e) + \psi_2 \log \det(G_W^e) + \log \det(G_\lambda^e) + \frac{1}{D} \psi_1 \psi_2 \lambda \sum (\eta_h^{\alpha(a)})^2 + \frac{1}{2D} \left( \lambda_i^{\alpha(a)} (G_\lambda^{e-1})_{ii'}^{\alpha\beta,ab} \lambda_{i'}^{\beta(b)} \right). \quad (69)$$

### 3.4.2 Expression of the action and the prefactor in terms of order parameters

This time, because of the two different systems, the order parameters carry an additional index  $a$ , which turns them into  $2 \times 2$  block matrices:

$$1 = \int dQ_{\alpha\beta}^{(ab)} d\hat{Q}_{\alpha\beta}^{(ab)} e^{\hat{Q}_{\alpha\beta}^{(ab)} (PQ_{\alpha\beta}^{(ab)} - \eta_h^{\alpha(a)} \eta_h^{\beta(b)})}, \quad (70)$$

$$1 = \int dR_{\alpha\beta}^{(ab)} d\hat{R}_{\alpha\beta}^{(ab)} e^{\hat{R}_{\alpha\beta}^{(ab)} (dR_{\alpha\beta}^{(ab)} - \lambda_i^{\alpha(a)} \lambda_i^{\beta(b)})}. \quad (71)$$

The systems being decoupled, we make the following ansatz for the order parameters:

$$Q = \left[ \begin{array}{ccc|ccc} q & \cdots & 0 & \tilde{q} & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & q & 0 & \cdots & \tilde{q} \\ \hline \tilde{q} & \cdots & 0 & q & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \tilde{q} & 0 & \cdots & q \end{array} \right], \quad R = \left[ \begin{array}{ccc|ccc} r & \cdots & 0 & \tilde{r} & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & r & 0 & \cdots & \tilde{r} \\ \hline \tilde{r} & \cdots & 0 & r & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & \tilde{r} & 0 & \cdots & r \end{array} \right]. \quad (72)$$

In virtue of the simple structure of the above matrices, the replica indices  $\alpha$  trivialize and we may replace the matrices  $Q$  and  $R$  by the  $2 \times 2$  matrices:

$$Q = \begin{bmatrix} q & \tilde{q} \\ \tilde{q} & q \end{bmatrix}, \quad R = \begin{bmatrix} r & \tilde{r} \\ \tilde{r} & r \end{bmatrix}.$$

Define:

$$[M_X^e]_{(ab)} \equiv \frac{1}{D} \lambda_i^{\alpha(a)} [G_X^{e-1}]_{ij} \lambda_j^{\beta(b)} = [R(I + \mu_1^2\psi_1 R)^{-1}]_{(ab)}, \quad (73)$$

$$[M_W^e]_{(ab)} \equiv \frac{1}{P} \eta^{\alpha(a)} [G_W^{e-1}]^{(ab)} \eta^{\beta(b)} = [Q(I + \psi_1 A^e Q)^{-1}]_{(ab)}, \quad (74)$$

where products are now over  $2 \times 2$  matrices. Then, one has:

$$P_{\Psi_3^e} [Q, R] := D \mu_1^2 \psi_1^2 \psi_2 R^{(12)} \left[ \mu_*^2 M_W^e{}^{(12)} + \mu_1^2 \left( M_X^e{}^{(12)} + \mu_1^2 \mu_*^2 \psi_1^2 [M_X^e M_W^e M_X^e]^{(12)} \right) + 2\mu_1^2 \mu_*^2 \psi_1 [M_X^e M_W^e]^{(12)} \right],$$

$$S^e [Q, R] := \psi_2 \log \det(G_X^e) + \psi_2 \log \det(G_W^e) + \sum_a \left[ (\psi_1^2 \psi_2 \lambda) \text{Tr} Q^{(a)(a)} + \text{Tr} \left( R^{(a)(a)} (Q^{-1})^{(a)(a)} \right) + \log \det Q^{(a)(a)} \right] - \psi_1 \log \det Q - \log \det R.$$

Where we have:

$$M_X^e = \frac{1}{(1 + \mu_1^2 \psi_1 r)^2 - (\mu_1^2 \psi_1 \tilde{r})^2} \begin{bmatrix} r + \mu_1^2 \psi_1 (r^2 - \tilde{r}^2) & \tilde{r} \\ \tilde{r} & r + \mu_1^2 \psi_1 (r^2 - \tilde{r}^2) \end{bmatrix}, \quad (75)$$

$$A_{(ab)}^e = \delta_{(ab)} - \mu_1^2 \psi_1 [M_X^e]_{(ab)}, \quad (76)$$

$$[M_W^e]_{(ab)} = q \delta_{(ab)} - \mu_\star^2 \psi_1 [A^e (I + \mu_\star^2 \psi_1 q A^e)^{-1}]_{(ab)}, \quad (77)$$

$$\det(G_X^e) = \det(\delta_{(ab)} + \psi_1 \mu_1^2 R_{(ab)}), \quad (78)$$

$$\det(G_W^e) = \det(\delta_{ab} + \psi_1 q A_{(ab)}^e). \quad (79)$$

Finally, we are left with:

$$S^e(q, r, \tilde{q}, \tilde{r}) = n (S_0(q, r) + S_1^e(q, r, \tilde{q}, \tilde{r})),$$

$$S_1^e(q, r, \tilde{q}, \tilde{r}) = \tilde{r}^2 f^e(q, r) + \tilde{q}^2 g^e(q, r),$$

$$f^e(q, r) = \frac{2r \mu_1^2 \psi_1 (1 + q \mu_\star^2 \psi_1) + (1 + q \mu_\star^2 \psi_1)^2 - r^2 \mu_1^4 \psi_1^2 (-1 + \psi_2)}{2r^2 (1 + r \mu_1^2 \psi_1 + q \mu_\star^2 \psi_1)^2},$$

$$g^e(q, r) = \frac{\psi_1}{2q^2}.$$

Where  $S_0$  was defined in (52).

### 3.4.3 Expression of the ensembling terms

Evaluating the fluctuations around the saddle point follows through in the same way as for the vanilla terms, with the following expressions of the prefactors:

$$P_{\Psi_2^e} [Q, R] = D \psi_1^2 \psi_2 \mu_1^2 \tilde{r} [\mu_1^2 P_{XX}^e - 2 \mu_1^2 \mu_\star^2 \psi_1 P_{WX}^e + \mu_\star^2 P_{WW}^e], \quad (80)$$

$$P_{\Psi_3^e} [Q, R] = D \mu_1^2 \psi_1^2 \psi_2 R^{(12)} \left[ \mu_\star^2 M_W^{e(12)} + \mu_1^2 \left( M_X^{e(12)} + \mu_1^2 \mu_\star^2 \psi_1^2 [M_X^e M_W^e M_X^e]^{(12)} \right) + 2 \mu_1^2 \mu_\star^2 \psi_1 [M_X^e M_W^e]^{(12)} \right], \quad (81)$$

$$P_{XX}^e = \psi_2 N_X^{e12} + M_X^{e12} + 2 \psi_2 (\mu_1 \mu_\star \psi_1)^2 [M_X^e N_X^e M_W^e]^{12} + (\mu_1 \mu_\star \psi_1)^2 [M_X^e M_W^e M_X^e]^{12} \\ + \psi_2 (\mu_1 \mu_\star \psi_1)^4 [M_X^e M_W^e N_X^e M_W^e M_X^e]^{12}, \quad (82)$$

$$P_{WX}^e = \psi_2 [N_X^e M_W^e]^{12} + [M_X^e M_W^e]^{12} + \psi_2 (\mu_1 \mu_\star \psi_1)^2 [M_X^e M_W^e N_X^e M_W^e]^{12}, \quad (83)$$

$$P_{WW}^e = [M_W^e]^{12} + \psi_2 (\mu_1 \mu_\star \psi_1)^2 [M_W^e N_X^e M_W^e]^{12}. \quad (84)$$

## 3.5 Computation of the *divide and conquer* term

Here, we are interested in computing the term  $\Psi_2^d$ . This term differs from the previous ones in that there are now two independent data matrices  $\mathbf{X}^{(1)}$  and  $\mathbf{X}^{(2)}$ . The calculations for the action and the prefactor are very similar to calculations performed for the *ensembling* terms  $\Psi_2^e, \Psi_3^e$ , with the addition that  $\mathbf{X}$  now also carries an index ( $a$ ).

Firstly let us write  $\Psi_2^d$  as a trace over random matrices:

$$\Psi_2^d = \frac{\mu_1^2}{d^2} \text{Tr} \left[ \mathbf{X}^{(1)} \mathbf{Z}^{(1)} \mathbf{B}^{(1)-1} \boldsymbol{\Theta}^{(1)} \boldsymbol{\Theta}^{(2)} \mathbf{B}^{(2)-1} \mathbf{Z}^{(2)} \mathbf{X}^{(2)} \right]. \quad (85)$$



Calculations follow through in the same way as in the previous sections. Using the replica formula (55), and performing the integrals over the Gaussian variables the following quantities appear:

$$(G_X^d)_{ii'}^{(ab)} = \delta_{ii'} + \frac{\mu_1^2 \psi_1}{D} \sum_{\alpha} \lambda_i^{\alpha(a)} \lambda_{i'}^{\alpha(a)}, \quad (86)$$

$$(J_X^d)_i^{(a)} = \frac{\mu_1 \mu_{\star} \sqrt{\psi_1}}{D^2} \sum_{\alpha} \lambda_i^{\alpha(a)} \eta_h^{\alpha(a)} W_h^{(a)}, \quad (87)$$

$$(G_W^d)_{hh'}^{(ab)} = \delta_{hh'} + \frac{\mu_{\star}^2}{D} \sum_{\alpha\beta} \eta_h^{\alpha(a)} A^{\alpha\beta, ab} \eta_{h'}^{\beta(b)}, \quad (88)$$

$$(J_W^d)_h = 0, \quad (89)$$

$$(G_{\Theta}^d)_{hh', ii'}^{(ab)} = \delta_{hh', ii'}^{ab}, \quad (90)$$

$$(J_{\Theta}^d)_{hi}^{(a)} = -i \sum_{\alpha} \hat{\lambda}_i^{\alpha(a)} \eta_h^{\alpha(a)}, \quad (91)$$

$$(G_{\lambda}^d)_{ii'}^{\alpha\beta, (ab)} = 2\delta^{ab} \delta^{ii'} \eta_h^{\alpha(a)} \eta_h^{\beta(b)}, \quad (92)$$

$$(J_{\lambda}^d)_i^{\alpha(a)} = i\sqrt{P} \lambda_i^{\alpha(a)}. \quad (93)$$

$$(94)$$

The saddle point ansatz for  $Q$  and  $R$  is the same as the one for the ensembling terms (see (72)). The procedure to evaluate  $\Psi_2^d$  is also the same as the one for  $\Psi_2^e$  except the Hessian is taken with respect to  $S^d$ . The final result is given below.

$$\begin{aligned} P_{\Psi_2^d}[Q, R] &= D\mu_1^2 \psi_1 \psi_2^2 \tilde{r} [\psi_1 \mu_1^2 P_{XX} + 2\mu_{\star}^2 \mu_1^2 \psi_1^2 P_{WX} + \mu_{\star}^2 \psi_1 P_{WW}], \\ P_{XX} &= \left( N_X^{11} + 2(\mu_1 \mu_{\star} \psi_1)^2 [N_X M_W M_X]^{11} + (\mu_1 \mu_{\star} \psi_1)^4 [M_X M_W N_X M_W M_X]^{11} \right), \\ P_{WX} &= [N_X M_W]^{11} + (\mu_1 \mu_{\star} \psi_1)^2 [M_X M_W N_X M_W]^{11}, \\ P_{WW} &= (\mu_1 \mu_{\star} \psi_1)^2 [M_W N_X M_W]^{11}, \\ S^d[q, r, \tilde{q}, \tilde{r}] &= n (S_0(q, r) + S_1^d(q, r, \tilde{q}, \tilde{r})), \\ S_1^d(q, r, \tilde{q}, \tilde{r}) &= \frac{\tilde{r}^2}{r^2} + \frac{\psi_1 \tilde{q}^2}{q^2}. \end{aligned}$$

With:

$$\begin{aligned} M_X &= \frac{r}{1 + \mu_1^2 \psi_1 r}, \\ A &= 1 - \mu_1^2 \psi_1 M_X, \\ M_W &= \frac{q}{1 + \mu_{\star}^2 \psi_1 q A}, \\ N_X &= \frac{\tilde{r}}{(1 + \mu_1^2 \psi_1 r)^2}. \end{aligned}$$

## References

- [1] Nitesh V Chawla, Thomas E Moore, Lawrence O Hall, Kevin W Bowyer, W Philip Kegelmeyer, and Clayton Springer. Distributed learning with bagging-like performance. *Pattern recognition letters*, 24(1-3):455–471, 2003.
- [2] Harris Drucker, Corinna Cortes, Lawrence D Jackel, Yann LeCun, and Vladimir Vapnik. Boosting and other ensemble methods. *Neural Computation*, 6(6):1289–1301, 1994.
- [3] Mario Geiger, Arthur Jacot, Stefano Spigler, Franck Gabriel, Levent Sagun, Stéphane d’Ascoli, Giulio Biroli, Clément Hongler, and Matthieu Wyart. Scaling description of generalization with number of parameters in deep learning. *arXiv preprint arXiv:1901.01608*, 2019.
- [4] Song Mei and Andrea Montanari. The generalization error of random features regression: Precise asymptotics and double descent curve. *arXiv preprint arXiv:1908.05355*, 2019.
- [5] Yuchen Zhang, John Duchi, and Martin Wainwright. Divide and conquer kernel ridge regression: A distributed algorithm with minimax optimal rates. *The Journal of Machine Learning Research*, 16(1):3299–3340, 2015.