

## A Definitions from Linear Analysis

We denote column vectors and matrices with small and capital bold letters, respectively, e.g.  $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_d]^\top \in \mathbb{R}^d$  and  $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$ . The singular values of a rectangular matrix  $\mathbf{A} \in \mathbb{R}^{n \times d}$  are non-negative and are denoted  $s_{\max}(\mathbf{A}) = s_1(\mathbf{A}) \geq \dots \geq s_{n \wedge d}(\mathbf{A}) = s_{\min}(\mathbf{A})$ . The rank of  $\mathbf{A}$  is  $r = \max\{k \mid s_k(\mathbf{A}) > 0\}$ . The eigenvalues of a Positive Semi-Definite (PSD) matrix  $\mathbf{M} \in \mathbb{R}^{d \times d}$  are non-negative and are denoted  $\lambda_{\max}(\mathbf{M}) = \lambda_1(\mathbf{M}) \geq \dots \geq \lambda_d(\mathbf{M}) = \lambda_{\min}(\mathbf{M})$ , while  $\lambda_{\min}^+(\mathbf{M})$  denotes the smallest *positive (non-zero)* eigenvalue.

When  $\mathbf{M} \in \mathbb{R}^{d \times d}$  is positive definite, we define  $\|\mathbf{x}\|_{\mathbf{M}}$  for  $\mathbf{x} \in \mathbb{R}^d$  by  $\|\mathbf{x}\|_{\mathbf{M}} = \sqrt{\mathbf{x}^\top \mathbf{M} \mathbf{x}}$ . It is easy to check that  $\|\cdot\|_{\mathbf{M}}$  is indeed a norm on  $\mathbb{R}^d$ , hence it induces a metric over  $\mathbb{R}^d$ , with the distance between  $\mathbf{x}$  and  $\mathbf{y}$  given by  $\|\mathbf{x} - \mathbf{y}\|_{\mathbf{M}} = \sqrt{(\mathbf{x} - \mathbf{y})^\top \mathbf{M} (\mathbf{x} - \mathbf{y})}$ . If  $\mathbf{M}$  is only semi-definite, these definitions would give a semi-norm and semi-metric. Note that  $\|\mathbf{x}\|_{\mathbf{M}} = \|\mathbf{M}^{1/2} \mathbf{x}\|$  where  $\mathbf{M}^{1/2}$  is the matrix square root of  $\mathbf{M}$ . If we set  $\mathbf{M} = \mathbf{I}$ , the identity matrix, then the norm  $\|\cdot\|_{\mathbf{M}}$  reduces to the standard Euclidean norm:  $\|\mathbf{x}\| = \sqrt{\mathbf{x}^\top \mathbf{x}}$ .

Combining the Cauchy-Schwarz inequality and the definition of operator norm  $\|\mathbf{M}\| = s_{\max}(\mathbf{M})$ , which implies  $\|\mathbf{M}\mathbf{x}\| \leq \|\mathbf{M}\| \|\mathbf{x}\|$ , we get the inequality  $\|\mathbf{x}\|_{\mathbf{M}}^2 \leq \|\mathbf{x}\|^2 \|\mathbf{M}\|$ .

## B Regression Model

Consider a linear regression model

$$Y = \mathbf{X}^\top \mathbf{w}^* + \varepsilon,$$

where  $\mathbf{w}^* \in \mathbb{R}^d$  is fixed and unknown, the random input  $\mathbf{X} \in \mathbb{R}^d$  is distributed according to some unknown distribution  $P_X$  supported on a unit ball, and noise is given by a random variable  $\varepsilon$  independent from  $\mathbf{X}$ , such that  $\mathbb{E}[\varepsilon] = 0$  and  $\mathbb{E}[\varepsilon^2] = \sigma^2$ . After observing an i.i.d. training sample  $S = ((\mathbf{X}_i, Y_i))_{i=1}^n$ , we run GD on the empirical squared loss:

$$\hat{L}_S(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{X}_i - Y_i)^2.$$

The sample covariance matrix is  $\hat{\boldsymbol{\Sigma}} = (\mathbf{X}_1 \mathbf{X}_1^\top + \dots + \mathbf{X}_n \mathbf{X}_n^\top) / n$ , and its population counterpart is the covariance matrix  $\boldsymbol{\Sigma} = \mathbb{E}[\mathbf{X} \mathbf{X}^\top]$ . Let  $\hat{\boldsymbol{\Sigma}} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$  be the SVD of  $\hat{\boldsymbol{\Sigma}}$  with orthogonal matrices  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_d]$  and  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_d]$ , and scaling matrix  $\mathbf{S} = \text{diag}(s_1(\hat{\boldsymbol{\Sigma}}), \dots, s_d(\hat{\boldsymbol{\Sigma}}))$  of singular values arranged in decreasing order:  $s_1(\hat{\boldsymbol{\Sigma}}) \geq s_2(\hat{\boldsymbol{\Sigma}}) \geq \dots \geq s_d(\hat{\boldsymbol{\Sigma}}) \geq 0$ . Note that  $s_i(\hat{\boldsymbol{\Sigma}}) = \lambda_i(\hat{\boldsymbol{\Sigma}}) =: \hat{\lambda}_i$  since  $\hat{\boldsymbol{\Sigma}}$  is positive semi-definite. The matrix  $\hat{\boldsymbol{\Sigma}}$  might be degenerate ( $\hat{\lambda}_d = 0$ ), and the non-degenerate part is given by  $\mathbf{U}_r := [\mathbf{u}_1, \dots, \mathbf{u}_r]$ ,  $\mathbf{V}_r := [\mathbf{v}_1, \dots, \mathbf{v}_r]$  and  $\mathbf{S}_r := \text{diag}(\hat{\lambda}_1, \dots, \hat{\lambda}_r)$ , where  $r = \text{rank}(\hat{\boldsymbol{\Sigma}})$ . We denote  $\hat{\mathbf{M}} = \mathbf{U}_r \mathbf{U}_r^\top$ , and note that  $\hat{\mathbf{M}}^2 = \hat{\mathbf{M}}$ . For the minimal *positive (non-zero)* eigenvalue we use the shorthand  $\hat{\lambda}_{\min}^+ = \lambda_{\min}^+(\hat{\boldsymbol{\Sigma}}) = \lambda_r(\hat{\boldsymbol{\Sigma}})$ .

## C Proof of Theorem 2 (Excess Risk of GD)

In this section we consider the standard GD algorithm, that is  $\mathcal{A}_S(\mathbf{w}_0) = \mathbf{w}_T$ , which is obtained recursively by applying the update rule  $\mathbf{w}_{t+1} = \mathbf{w}_t - \alpha \nabla \hat{L}_S(\mathbf{w}_t)$  with some step size  $\alpha > 0$  and initialization  $\mathbf{w}_0 \in \mathbb{R}^d$ . The rule is iterated for  $t = 0, \dots, T-1$ .

Recall that the *excess risk* of  $\mathcal{A}_S(\mathbf{w}_0)$  is defined as

$$\mathcal{E}_T = L(\mathcal{A}_S(\mathbf{w}_0)) - L(\mathbf{w}^*).$$

Next we give upper bounds on the excess risk of GD output.

**Theorem 2 (restated).** Assume that  $\alpha \leq 1/\widehat{\lambda}_1$ . For any  $\mathbf{w}_0 \in \mathbb{R}^d$  and  $x > 0$ , with probability at least  $1 - e^{-x}$  over inputs we have

$$\begin{aligned} \mathbb{E}_\epsilon[\mathcal{E}_T] &\leq \widehat{\lambda}_1^2(1 - \alpha\widehat{\lambda}_{\min}^+)^{2T} \|\mathbf{w}^* - \mathbf{w}_0\|^2 + \frac{4\sigma^2}{n} \left( \frac{\widehat{\lambda}_1}{\widehat{\lambda}_{\min}^+} \right)^2 \\ &\quad + \frac{12}{\sqrt{n}} \left( \sqrt{\ln d} + \sqrt{x} \right) \left( \|\mathbf{w}_0\|^2 + 2\|\mathbf{w}^*\|^2 + \sigma^2 \sum_{i=1}^r \widehat{\lambda}_i^{-1} \right). \end{aligned}$$

The proof of Theorem 2 is based on a the following decomposition, and remaining subsections will be dedicated to bounding the constituent terms.

**Proposition 1 (restated).** For any  $\mathbf{w}_0 \in \mathbb{R}^d$ ,

$$\mathcal{E}_T \leq \underbrace{\|\mathcal{A}_S(\mathbf{w}_0) - \mathcal{A}_S(\mathbf{w}^*)\|_{\widehat{\Sigma}}^2}_{(1)} + \underbrace{\|\mathcal{A}_S(\mathbf{w}^*) - \mathbf{w}^*\|_{\widehat{\Sigma}}^2}_{(2)} + \underbrace{\|\Sigma - \widehat{\Sigma}\|_2 (\|\mathcal{A}_S(\mathbf{w}_0)\|^2 + \|\mathbf{w}^*\|^2)}_{(3)}.$$

*Proof.* Observe that for the square loss we have

$$\begin{aligned} L(\mathcal{A}_S(\mathbf{w}_0)) - L(\mathbf{w}^*) &= \frac{1}{2} \|\mathcal{A}_S(\mathbf{w}_0) - \mathbf{w}^*\|_{\Sigma}^2 \\ &= \frac{1}{2} \|\mathcal{A}_S(\mathbf{w}_0) - \mathbf{w}^*\|_{\widehat{\Sigma}}^2 + \frac{1}{2} \|\mathcal{A}_S(\mathbf{w}_0) - \mathbf{w}^*\|_{\Sigma - \widehat{\Sigma}}^2 \\ &= \frac{1}{2} \underbrace{\|\mathcal{A}_S(\mathbf{w}_0) - \mathbf{w}^*\|_{\widehat{\Sigma}}^2}_{(a)} + \frac{1}{2} \underbrace{\|\mathcal{A}_S(\mathbf{w}_0) - \mathbf{w}^*\|_{\Sigma - \widehat{\Sigma}}^2}_{(b)}. \end{aligned}$$

Note that bounding term (a) reduces to

$$\|\mathcal{A}_S(\mathbf{w}_0) - \mathbf{w}^*\|_{\widehat{M}}^2 \leq 2\|\mathcal{A}_S(\mathbf{w}_0) - \mathcal{A}_S(\mathbf{w}^*)\|_{\widehat{\Sigma}}^2 + 2\|\mathcal{A}_S(\mathbf{w}^*) - \mathbf{w}^*\|_{\widehat{\Sigma}}^2.$$

On the other we handle term (b) by Cauchy-Schwarz and triangle inequalities:

$$\|\mathcal{A}_S(\mathbf{w}_0) - \mathbf{w}^*\|_{\Sigma - \widehat{\Sigma}}^2 \leq 2\|\Sigma - \widehat{\Sigma}\| (\|\mathcal{A}_S(\mathbf{w}_0)\|^2 + \|\mathbf{w}^*\|^2).$$

□

Here, the first term (1), vanishes as long as the algorithm-map  $\mathcal{A}_S$  is contractive on the aforementioned subspace, which we show in Appendix C.1 thanks to the closed-form expression of GD iterates.

Term (2) captures algorithm's sensitivity to the label noise: How far would GD go when initialized at the *minimum of the risk*? Indeed it is easy to see that when there is no label noise, term (2) is zero (one can demonstrate it by the descent lemma). In the presence of noise, matters are more complicated, and we employ more or less standard technique where we recursively track the distance between GD iterates and "virtual" iterates, which are obtained as if we could remove the noise. This results in a bound  $4(\sigma^2/n)(\widehat{\lambda}_{\min}^+)^{-2}$ , which is shown in Appendix C.3.

Finally, term (3) is essentially a concentration of the sample covariance matrix and ensuring that the norm of the solution remains well-behaved, see Appendix C.2. The concentration of the sample covariance matrix is due to the matrix Chernoff inequality [Tropp, 2012]. We control the norm of the solution by relating it to the Moore-Penrose pseudoinverse solution which can be written in a closed-form. This makes it easy to see that when the label noise is absent, the norm depends only on  $\mathbf{w}^*$  and  $\mathbf{w}_0$ . On the other hand, when the noise is present things will depend on behaviour of  $\widehat{\lambda}_{\min}^+$  as discussed in Section 3.2.

### C.1 Contractivity of GD (Term (1))

Contractivity of GD comes from the following straightforward proposition.

**Proposition 2.** For a  $T$ -step gradient descent map  $\mathcal{A}_S : \mathbb{R}^d \rightarrow \mathbb{R}^d$  with step size  $\alpha > 0$  applied to the least squares, and for all  $\mathbf{w}_0 \in \mathbb{R}^d$ , we have a.s. that

$$\mathcal{A}_S(\mathbf{w}_0) = (\mathbf{I} - \alpha \widehat{\Sigma})^T \mathbf{w}_0 + \alpha \sum_{t=0}^{T-1} (\mathbf{I} - \alpha \widehat{\Sigma})^t \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i Y_i \right).$$

*Proof.* Abbreviate  $\mathbf{C} = (\mathbf{X}_1 Y_1 + \dots + \mathbf{X}_n Y_n)/n$ . Since  $\nabla \widehat{L}_S(\mathbf{w}) = \widehat{\Sigma} \mathbf{w} - \mathbf{C}$ , observe that

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \alpha (\widehat{\Sigma} \mathbf{w}_{t-1} - \mathbf{C}) = (\mathbf{I} - \alpha \widehat{\Sigma}) \mathbf{w}_{t-1} + \alpha \mathbf{C}.$$

A simple recursive argument reveals that for every  $\mathbf{w}_0 \in \mathbb{R}^d$

$$\begin{aligned} \mathcal{A}_S(\mathbf{w}_0) &= \mathbf{w}_T = (\mathbf{I} - \alpha \widehat{\Sigma}) \mathbf{w}_{T-1} + \alpha \mathbf{C} \\ &= (\mathbf{I} - \alpha \widehat{\Sigma})^2 \mathbf{w}_{T-2} + \alpha (\mathbf{I} - \alpha \widehat{\Sigma}) \mathbf{C} + \alpha \mathbf{C} \\ &= (\mathbf{I} - \alpha \widehat{\Sigma})^3 \mathbf{w}_{T-3} + \alpha (\mathbf{I} - \alpha \widehat{\Sigma})^2 \mathbf{C} + \alpha (\mathbf{I} - \alpha \widehat{\Sigma}) \mathbf{C} + \alpha \mathbf{C} \\ &\dots \\ &= (\mathbf{I} - \alpha \widehat{\Sigma})^T \mathbf{w}_0 + \alpha \sum_{t=0}^{T-1} (\mathbf{I} - \alpha \widehat{\Sigma})^t \mathbf{C}. \end{aligned}$$

□

Proposition 2 implies the following simple fact.

**Corollary 1** (Contractivity of GD). *The  $T$ -step gradient descent map  $\mathcal{A}_S : \mathbb{R}^d \rightarrow \mathbb{R}^d$  with step size  $\alpha > 0$  applied to the least squares problem satisfies, for all  $\mathbf{w}_0, \mathbf{u}_0 \in \mathbb{R}^d$ ,*

$$\|\mathcal{A}_S(\mathbf{w}_0) - \mathcal{A}_S(\mathbf{u}_0)\|_{\widehat{\Sigma}} \leq \widehat{\lambda}_1 (1 - \alpha \widehat{\lambda}_{\min}^+)^T \|\mathbf{w}_0 - \mathbf{u}_0\|.$$

*Proof.* Clearly  $\|\mathcal{A}_S(\mathbf{w}_0) - \mathcal{A}_S(\mathbf{u}_0)\|_{\widehat{\Sigma}} \leq \widehat{\lambda}_1 \|\mathcal{A}_S(\mathbf{w}_0) - \mathcal{A}_S(\mathbf{u}_0)\|_{\mathbf{U}_r \mathbf{U}_r^\top}$ . By Proposition 2 for any  $\mathbf{w}_0, \mathbf{u}_0 \in \mathbb{R}^d$ :

$$\begin{aligned} \|\mathcal{A}_S(\mathbf{w}_0) - \mathcal{A}_S(\mathbf{u}_0)\|_{\mathbf{U}_r \mathbf{U}_r^\top} &= \|(\mathbf{I} - \alpha \widehat{\Sigma})^T (\mathbf{w}_0 - \mathbf{u}_0)\|_{\mathbf{U}_r \mathbf{U}_r^\top} \\ &= \|\mathbf{U}_r^\top (\mathbf{I} - \alpha \widehat{\Sigma})^T (\mathbf{w}_0 - \mathbf{u}_0)\| \\ &\leq \|\mathbf{U}_r^\top (\mathbf{I} - \alpha \widehat{\Sigma})^T\| \|\mathbf{w}_0 - \mathbf{u}_0\|. \end{aligned}$$

Now,

$$\mathbf{U}_r^\top (\mathbf{I} - \alpha \widehat{\Sigma})^T = \mathbf{U}_r^\top \mathbf{U} (\mathbf{I} - \alpha \mathbf{S})^T \mathbf{V}^\top = \mathbf{I}_{r \times d} (\mathbf{I} - \alpha \mathbf{S})^T \mathbf{V}^\top = (\mathbf{I}_{r \times d} - \alpha \mathbf{S}_{r \times d})^T \mathbf{V}^\top$$

where subscript  $r \times n$  stands for clipping the matrix to  $r$  rows and  $d$  columns. The above implies that the operator norm of  $\mathbf{U}_r^\top (\mathbf{I} - \alpha \widehat{\Sigma})^T$  satisfies  $\|\mathbf{U}_r^\top (\mathbf{I} - \alpha \widehat{\Sigma})^T\|_2 \leq (1 - \alpha \lambda_{\min}^+(\widehat{\Sigma}))^T$ . □

## C.2 Concentration of Spectrum (Term (3))

Our next goal is to understand the behaviour of

$$\|\Sigma - \widehat{\Sigma}\|_2 (\|\mathcal{A}_S(\mathbf{w}_0)\|^2 + \|\mathbf{w}^*\|^2). \quad (2)$$

A high-probability concentration of a covariance matrix is readily given by the matrix Chernoff inequality:

**Theorem 3** (Tropp [2012]). *Suppose that inputs have a bounded spectral norm which is at most  $B_X$ . Then,*

$$\mathbb{P} \left( \|\widehat{\Sigma} - \Sigma\|_2 \geq \frac{6B_X}{\sqrt{n}} (\sqrt{\ln d} + \sqrt{x}) \right) \leq e^{-x}.$$

Next we need to control  $\|\mathcal{A}_S(\mathbf{w}_0)\|^2$ . This will be done by using a closed-form of the GD iterate at step  $T$  given by Proposition 2:

$$\|\mathcal{A}_S(\mathbf{w}_0)\|^2 \leq 2 \underbrace{\|(\mathbf{I} - \alpha \widehat{\Sigma})^T \mathbf{w}_0\|_2^2}_{(a)} + 2 \underbrace{\left\| \alpha \sum_{t=0}^{T-1} (\mathbf{I} - \alpha \widehat{\Sigma})^t \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i Y_i \right) \right\|_2^2}_{(b)}$$

where Cauchy-Schwarz inequality gives (a)  $\leq \|\mathbf{w}_0\|^2$ , while (b) converges to the Moore-Penrose pseudoinverse of  $\widehat{\Sigma}$  as  $T \rightarrow \infty$  as long as  $\alpha \leq 1/\widehat{\lambda}_1$  [Ben-Israel and Charnes, 1963]. Hence, we will need to quantify its squared  $\ell_2$  norm:

**Proposition 3.** *Suppose that given some  $\mathbf{w}^* \in \mathbb{R}^d$ ,  $\mathbf{D}^\top = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ , and  $\boldsymbol{\varepsilon} \in \mathbb{R}^n$ , we have  $\mathbf{y} = \mathbf{D}\mathbf{w}^* + \boldsymbol{\varepsilon}$ . Then, the Moore-Penrose pseudoinverse solution  $\mathbf{w}^{\text{pinv}} = (\mathbf{D}^\top \mathbf{D})^\dagger \mathbf{D}^\top \mathbf{y}$  satisfies*

$$\|\mathbf{w}^{\text{pinv}}\|^2 = \|\mathbf{w}^*\|^2 + 2\boldsymbol{\varepsilon}^\top (\mathbf{D}^\top \mathbf{D})^\dagger \mathbf{w}^* + \|\boldsymbol{\varepsilon}\|_{(\mathbf{D}\mathbf{D}^\top)^{-1}}^2.$$

*Proof.* Observe that

$$\begin{aligned} \|\mathbf{w}^{\text{pinv}}\|^2 &= \mathbf{y}^\top \mathbf{D} (\mathbf{D}^\top \mathbf{D})^\dagger \mathbf{D}^\top \mathbf{y} \\ &= (\mathbf{D}\mathbf{w}^* + \boldsymbol{\varepsilon})^\top \mathbf{D} (\mathbf{D}^\top \mathbf{D})^\dagger \mathbf{D}^\top (\mathbf{D}\mathbf{w}^* + \boldsymbol{\varepsilon}) \\ &= \|\mathbf{w}^*\|^2 + 2\boldsymbol{\varepsilon}^\top (\mathbf{D}^\top \mathbf{D})^\dagger \mathbf{D}^\top \mathbf{D}\mathbf{w}^* + \boldsymbol{\varepsilon}^\top \mathbf{D} (\mathbf{D}^\top \mathbf{D})^\dagger \mathbf{D}^\top \boldsymbol{\varepsilon} \\ &= \|\mathbf{w}^*\|^2 + 2\boldsymbol{\varepsilon}^\top (\mathbf{D}^\top \mathbf{D})^\dagger \mathbf{w}^* + \boldsymbol{\varepsilon}^\top (\mathbf{D}\mathbf{D}^\top)^{-1} \boldsymbol{\varepsilon}. \end{aligned}$$

□

Putting things together we have that w.p. at least  $1 - e^{-x}$  over  $(\mathbf{X}_1, \dots, \mathbf{X}_n)$ ,

$$\text{Eq. (2)} \leq \frac{12}{\sqrt{n}} \left( \sqrt{\ln d} + \sqrt{x} \right) \left( \|\mathbf{w}_0\|^2 + 2\|\mathbf{w}^*\|^2 + \frac{2}{n} \boldsymbol{\varepsilon}^\top \widehat{\Sigma}^\dagger \mathbf{w}^* + \|\boldsymbol{\varepsilon}\|_{(\mathbf{D}\mathbf{D}^\top)^{-1}}^2 \right).$$

Moreover, note that taking expectation over label noise gives

$$\mathbb{E}_\varepsilon[\text{Eq. (2)}] \leq \frac{12}{\sqrt{n}} \left( \sqrt{\ln d} + \sqrt{x} \right) \left( \|\mathbf{w}_0\|^2 + 2\|\mathbf{w}^*\|^2 + \sigma^2 \sum_{i=1}^r \widehat{\lambda}_i^{-1} \right).$$

### C.3 Bounding term (2)

The goal of this section is to bound

$$\|\mathcal{A}_S(\mathbf{w}^*) - \mathbf{w}^*\|_{\widehat{\Sigma}}^2,$$

which captures how sensitive the algorithm is to the noise when initialize at the optimum. When there is no label noise, the standard descent lemma (not shown here) readily gives that the term vanishes.

**Lemma 2** (Descent Lemma). *Assuming that  $\alpha \leq 1/\widehat{\lambda}_1$ ,*

$$\sum_{t=0}^{T-1} \|\nabla \widehat{L}_S(\mathbf{w}_t)\|^2 \leq \frac{2}{\alpha} \left( \widehat{L}_S(\mathbf{w}_0) - \widehat{L}_S(\mathcal{A}_S(\mathbf{w}_0)) \right).$$

In particular, if  $\mathbf{w}_t^*$  are the iterates of GD when starting from  $\mathbf{w}^*$  (so that  $\mathbf{w}_0 = \mathbf{w}^*$ ), then

$$\begin{aligned} \|\mathcal{A}_S(\mathbf{w}^*) - \mathbf{w}^*\|_M^2 &= \left\| \alpha \sum_{t=0}^{T-1} \nabla \widehat{L}_S(\mathbf{w}_t^*) \right\|_M^2 \\ &\leq \alpha^2 T \sum_{t=0}^{T-1} \|\nabla \widehat{L}_S(\mathbf{w}_t^*)\|_M^2 \\ &\leq \alpha^2 \frac{2T}{\alpha} \left( \widehat{L}_S(\mathbf{w}^*) - \widehat{L}_S(\mathbf{w}_T^*) \right) \\ &\leq 2\alpha T \widehat{L}_S(\mathbf{w}^*) = 0. \end{aligned}$$

Here we consider the case with i.i.d. label noise  $(\varepsilon_1, \dots, \varepsilon_n)$  and therefore show a more general handling of the term. Recall that  $\mathbb{E}[\varepsilon_i] = 0$  and  $\mathbb{E}[\varepsilon_i^2] = \sigma^2$  for  $i \in [n]$ . Throughout this section abbreviate  $\mathbb{E}[\cdot \mid \mathbf{X}_1, \dots, \mathbf{X}_n] = \mathbb{E}_\varepsilon[\cdot]$ .

**Lemma 3.** Let  $\widehat{\mathbf{M}}$  be defined as in Section 3.1. For any  $T > 0$ , GD achieves

$$\mathbb{E}_\varepsilon \left[ \|\mathbf{w}^* - \mathcal{A}_S(\mathbf{w}^*)\|_{\widehat{\Sigma}}^2 \right] \leq \frac{4\sigma^2}{n} \left( \frac{\widehat{\lambda}_1}{\widehat{\lambda}_{\min}^+} \right)^2.$$

*Proof.* We begin by noting that the integral form of Taylor theorem gives us that for any  $\mathbf{w}^* \in \arg \min_{\mathbf{w} \in \mathbb{R}^d} \widehat{L}(\mathbf{w})$  and any  $\mathbf{w} \in \mathbb{R}^d$ ,

$$\begin{aligned} \widehat{L}(\mathbf{w}) - \widehat{L}(\mathbf{w}^*) &= \frac{1}{2} (\mathbf{w} - \mathbf{w}^*)^\top \left( \int_0^1 \nabla^2 \widehat{L}(\tau \mathbf{w} + (1 - \tau) \mathbf{w}^*) d\tau \right) (\mathbf{w} - \mathbf{w}^*) \\ &\geq \frac{1}{2} \cdot \widehat{\lambda}_{\min}^+ (\mathbf{w} - \mathbf{w}^*)^\top \widehat{\mathbf{M}} (\mathbf{w} - \mathbf{w}^*). \end{aligned}$$

Thus, taking  $\mathbf{w} = \mathcal{A}_S(\mathbf{w}^*)$ , we have

$$\begin{aligned} \mathbb{E}_\varepsilon \left[ \|\mathbf{w}^* - \mathcal{A}_S(\mathbf{w}^*)\|_{\widehat{\mathbf{M}}}^2 \right] &\leq \frac{1}{\widehat{\lambda}_{\min}^+} \left( \mathbb{E}_\varepsilon \widehat{L}(\mathbf{w}^*) - \mathbb{E}_\varepsilon \left[ \widehat{L}(\mathcal{A}_S(\mathbf{w}^*)) \right] \right) \\ &= \frac{1}{\widehat{\lambda}_{\min}^+} \left( \sigma^2 - \mathbb{E}_\varepsilon \left[ \widehat{L}(\mathcal{A}_S(\mathbf{w}^*)) \right] \right). \end{aligned}$$

Now, let's focus on the loss term on the r.h.s.:

$$\begin{aligned} \mathbb{E}_\varepsilon \left[ \widehat{L}_S(\mathbf{w}_T^*) \right] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_\varepsilon \left[ \left( (\mathbf{w}_T^* - \mathbf{w}_0^*)^\top \mathbf{X}_i - \varepsilon_i \right)^2 \right] \\ &= \sigma^2 - \frac{2}{n} \sum_{i=1}^n \mathbb{E}_\varepsilon \left[ \varepsilon_i (\mathbf{w}_T^* - \mathbf{w}_0^*)^\top \mathbf{X}_i \right] + \mathbb{E}_\varepsilon \left[ (\mathbf{w}_T^* - \mathbf{w}_0^*)^\top \widehat{\Sigma} (\mathbf{w}_T^* - \mathbf{w}_0^*) \right] \\ &\geq \sigma^2 - \frac{2}{n} \sum_{i=1}^n \mathbb{E}_\varepsilon \left[ \varepsilon_i (\mathbf{w}_T^* - \mathbf{w}_0^*)^\top \mathbf{X}_i \right] \\ &= \sigma^2 - \frac{2}{n} \sum_{i=1}^n \mathbb{E}_\varepsilon \left[ \varepsilon_i \mathbf{w}_T^{*\top} \mathbf{X}_i \right] \end{aligned}$$

where the last term is small when label noise is not too correlated with the output  $\mathbf{w}_T^*$ . Hence to control the term, we need to measure the effect of the noise on GD. To do so we will introduce an additional iterates  $(\tilde{\mathbf{w}}_t)_t$  constructed by running GD on labels without noise, that is

$$\tilde{\mathbf{w}}_{t+1}^* = \tilde{\mathbf{w}}_t^* - \alpha \nabla \tilde{L}_S(\tilde{\mathbf{w}}_t^*) \quad \text{where} \quad \tilde{L}(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^n \left( \mathbf{w}^\top \mathbf{X}_i - \mathbf{w}^{*\top} \mathbf{X}_i \right)^2.$$

The plan is then to bound the deviation  $\|\mathbf{w}_T^* - \tilde{\mathbf{w}}_T^*\|_{\widehat{\mathbf{M}}}$  which we will do recursively. We proceed:

$$\begin{aligned} &\frac{2}{n} \sum_{i=1}^n \mathbb{E}_\varepsilon \left[ \varepsilon_i \mathbf{w}_T^{*\top} \mathbf{X}_i \right] \\ &= \frac{2}{n} \sum_{i=1}^n \mathbb{E}_\varepsilon \left[ \varepsilon_i (\mathbf{w}_T^* - \tilde{\mathbf{w}}_T^*)^\top \mathbf{X}_i \right] \quad (\text{Note that } \mathbb{E}_\varepsilon[\tilde{\mathbf{w}}_T^* \mid \mathbf{X}_i] = 0) \\ &= \frac{2}{n} \sum_{i=1}^n \mathbb{E}_\varepsilon \left[ \varepsilon_i (\mathbf{w}_T^* - \tilde{\mathbf{w}}_T^*)^\top \widehat{\mathbf{M}} \mathbf{X}_i \right] \quad (\text{Since } \widehat{\mathbf{M}} \mathbf{X}_i = \mathbf{X}_i) \\ &\leq \frac{2}{n} \mathbb{E}_\varepsilon \left[ \left\| \sum_{i=1}^n \varepsilon_i \mathbf{X}_i \right\| \left\| \widehat{\mathbf{M}} (\mathbf{w}_T^* - \tilde{\mathbf{w}}_T^*) \right\| \right] \quad (\text{Cauchy-Schwarz}) \end{aligned}$$

Now we will handle  $\left\| \widehat{\mathbf{M}}(\mathbf{w}_T^* - \tilde{\mathbf{w}}_T^*) \right\| = \|\mathbf{w}_T^* - \tilde{\mathbf{w}}_T^*\|_{\widehat{\mathbf{M}}}$  by following a recursive argument. First, observe that for any  $t = 0, 1, 2, \dots$

$$\nabla \hat{L}(\tilde{\mathbf{w}}_t^*) = \widehat{\Sigma} \tilde{\mathbf{w}}_t^* - \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \mathbf{w}_0^* = \widehat{\Sigma}(\tilde{\mathbf{w}}_t^* - \mathbf{w}_0^*),$$

and at the same time

$$\nabla \hat{L}(\mathbf{w}_t^*) = \widehat{\Sigma} \mathbf{w}_t^* - \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top \mathbf{w}_0^* - \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \varepsilon_i = \widehat{\Sigma}(\mathbf{w}_t^* - \mathbf{w}_0^*) - \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \varepsilon_i.$$

Thus,

$$\|\mathbf{w}_{t+1}^* - \tilde{\mathbf{w}}_{t+1}^*\|_{\widehat{\mathbf{M}}} = \left\| \mathbf{w}_t^* - \tilde{\mathbf{w}}_t^* - \alpha \left( \nabla \hat{L}(\mathbf{w}_t^*) - \nabla \hat{L}(\tilde{\mathbf{w}}_t^*) \right) \right\|_{\widehat{\mathbf{M}}} \quad (3)$$

$$\begin{aligned} &= \left\| \mathbf{w}_t^* - \tilde{\mathbf{w}}_t^* - \alpha \widehat{\Sigma}(\mathbf{w}_t^* - \tilde{\mathbf{w}}_t^*) - \frac{\alpha}{n} \sum_{i=1}^n \mathbf{X}_i \varepsilon_i \right\|_{\widehat{\mathbf{M}}} \\ &= \left\| (\mathbf{I} - \alpha \widehat{\Sigma})(\mathbf{w}_t^* - \tilde{\mathbf{w}}_t^*) \right\|_{\widehat{\mathbf{M}}} + \frac{\alpha}{n} \left\| \sum_{i=1}^n \mathbf{X}_i \varepsilon_i \right\|_{\widehat{\mathbf{M}}} \\ &\stackrel{(a)}{\leq} \|\mathbf{I} - \alpha \widehat{\Sigma}\|_{\widehat{\mathbf{M}}} \|\mathbf{w}_t^* - \tilde{\mathbf{w}}_t^*\|_{\widehat{\mathbf{M}}} + \frac{\alpha}{n} \left\| \sum_{i=1}^n \mathbf{X}_i \varepsilon_i \right\|_{\widehat{\mathbf{M}}} \\ &\leq (1 - \alpha \widehat{\lambda}_{\min}^+) \|\mathbf{w}_t^* - \tilde{\mathbf{w}}_t^*\|_{\widehat{\mathbf{M}}} + \frac{\alpha}{n} \left\| \sum_{i=1}^n \mathbf{X}_i \varepsilon_i \right\|_{\widehat{\mathbf{M}}}. \end{aligned} \quad (4)$$

where in the step (a) we note that  $\widehat{\mathbf{M}}(\mathbf{I} - \alpha \widehat{\Sigma})(\mathbf{w}_t^* - \tilde{\mathbf{w}}_t^*) = \widehat{\mathbf{M}}(\mathbf{I} - \alpha \widehat{\Sigma})\widehat{\mathbf{M}}(\mathbf{w}_t^* - \tilde{\mathbf{w}}_t^*)$  (since  $\widehat{\mathbf{M}}^2 = \widehat{\mathbf{M}}$  and  $\widehat{\Sigma}\widehat{\mathbf{M}} = \widehat{\Sigma}$ ).

Now we use the fact that an elementary recursive relation  $x_{t+1} \leq a_t x_t + b_t$  with  $x_0 = 0$  unwinds to  $x_T \leq \sum_{t=1}^T b_t \prod_{k=t+1}^T a_k$ , which gives

$$\begin{aligned} \|\mathbf{w}_T^* - \tilde{\mathbf{w}}_T^*\|_{\widehat{\mathbf{M}}} &\leq \frac{\alpha}{n} \left\| \sum_{i=1}^n \mathbf{X}_i \varepsilon_i \right\|_{\widehat{\mathbf{M}}} \sum_{t=1}^T (1 - \alpha \widehat{\lambda}_{\min}^+)^{T-t} \\ &\leq \frac{\alpha}{n} \left\| \sum_{i=1}^n \mathbf{X}_i \varepsilon_i \right\|_{\widehat{\mathbf{M}}} \frac{1 - (1 - \alpha \widehat{\lambda}_{\min}^+)^T}{\alpha \widehat{\lambda}_{\min}^+}. \end{aligned}$$

Thus,

$$\begin{aligned} \frac{2}{n} \sum_{i=1}^n \mathbb{E}_{\varepsilon} \left[ \varepsilon_i (\mathbf{w}_T^* - \mathbf{w}_0^*)^\top \mathbf{X}_i \right] &\leq \frac{2}{n} \cdot \frac{1}{n} \mathbb{E}_{\varepsilon} \left[ \left\| \sum_{i=1}^n \mathbf{X}_i \varepsilon_i \right\|^2 \frac{1}{\widehat{\lambda}_{\min}^+} \right] \\ &\leq \frac{2\sigma^2}{n} \cdot \frac{1}{\widehat{\lambda}_{\min}^+} \end{aligned}$$

where we used a basic fact that

$$\mathbb{E}_{\varepsilon} \left[ \left\| \sum_{i=1}^n \mathbf{X}_i \varepsilon_i \right\|^2 \middle| \mathbf{X}_1, \dots, \mathbf{X}_n \right] = \sigma^2 \sum_{i=1}^n \|\mathbf{X}_i\|^2 \leq \sigma^2 n.$$

Putting all together completes the proof.  $\square$

## D Concentration of the Smallest Non-zero Eigenvalue: Proof

Here we show a non-asymptotic concentration of the smallest positive eigenvalue as stated in Section 3.3.

**Lemma 1 (restated).** Let  $\mathbf{D}^\top = [\mathbf{X}_1, \dots, \mathbf{X}_n] \in \mathbb{R}^{d \times n}$  be a matrix with i.i.d. columns, such that  $\max_i \|\mathbf{X}_i\|_{\psi_2} \leq K$ , and let  $\widehat{\Sigma} = \mathbf{D}^\top \mathbf{D} / n$ , and  $\Sigma = \mathbb{E}[\mathbf{X}_1 \mathbf{X}_1^\top]$ . Then, for every  $x \geq 0$ , with probability at least  $1 - 2e^{-x}$ , we have

$$\lambda_{\min}^+(\widehat{\Sigma}) \geq \lambda_{\min}^+(\Sigma) \left( 1 - K^2 \left( c\sqrt{\frac{d}{n}} + \sqrt{\frac{x}{n}} \right) \right)_+^2 \quad \text{for } n \geq d,$$

and furthermore, assuming that  $\|\mathbf{X}_i\|_{\Sigma^\dagger} = \sqrt{d}$  a.s. for all  $i \in [n]$ , we have

$$\lambda_{\min}^+(\widehat{\Sigma}) \geq \lambda_{\min}^+(\Sigma) \left( \sqrt{\frac{d}{n}} - K^2 \left( c + 6\sqrt{\frac{x}{n}} \right) \right)_+^2 \quad \text{for } n < d,$$

where we have an absolute constant  $c = 2^{3.5} \sqrt{\ln(9)}$ .

Lemma 1 proven shortly, is based upon the the next theorem gives us a non-asymptotic version of Bai-Yin law [Bai and Yin, 1993] for rectangular matrices whose rows are sub-Gaussian isotropic random vectors.

**Theorem 4** ([Vershynin, 2012, Theorem 5.39]). Let  $\mathbf{A} \in \mathbb{R}^{n \times d}$  whose rows  $(\mathbf{A}^\top)_i$  are independent sub-Gaussian isotropic random vectors in  $\mathbb{R}^d$ , such that  $K = \max_{i \in [n]} \|(\mathbf{A}^\top)_i\|_{\psi_2}$ . Then for every  $x \geq 0$ , with probability at least  $1 - 2e^{-x}$  one has

$$\sqrt{n} - 2^{3.5} K^2 (\sqrt{\ln(9)d} + \sqrt{x}) \leq s_{\min}(\mathbf{A}) \leq s_{\max}(\mathbf{A}) \leq \sqrt{n} + 2^{3.5} K^2 (\sqrt{\ln(9)d} + \sqrt{x}).$$

**Theorem 5** ([Vershynin, 2012, Theorem 5.58]). Let  $\mathbf{A} \in \mathbb{R}^{d \times n}$  whose columns  $\mathbf{A}_i$  are independent sub-Gaussian isotropic random vectors in  $\mathbb{R}^d$  with  $\|\mathbf{A}_i\| = \sqrt{d}$  a.s., such that  $K = \max_{i \in [n]} \|\mathbf{A}_i\|_{\psi_2}$ . Then for every  $x \geq 0$ , with probability at least  $1 - 2e^{-x}$  one has

$$\sqrt{d} - 2^{3.5} K^2 (\sqrt{\ln(9)n} + 6\sqrt{x}) \leq s_{\min}(\mathbf{A}) \leq s_{\max}(\mathbf{A}) \leq \sqrt{d} + 2^{3.5} K^2 (\sqrt{\ln(9)n} + 6\sqrt{x})$$

Above two theorems lead to the following non-asymptotic version of a Bai-Yin law.

*Proof of Lemma 1.* The proof considers two cases: 1) when number of observations exceeds the dimension, which is handled by the concentration of a minimal non-zero eigenvalue of a covariance matrix; 2) when dimension exceeds number of observations, which is handled by concentration of the Gram matrix.

**Case  $n \geq d$ .** We will apply Theorem 4 with  $\mathbf{A} = (\Sigma^{\dagger \frac{1}{2}} \mathbf{D}^\top)^\top = \mathbf{D} \Sigma^{\dagger \frac{1}{2}}$  whose rows are independent and isotropic, and in addition by Cauchy-Schwarz inequality:

$$\|\Sigma^{\dagger \frac{1}{2}}\| s_{\min}(\mathbf{D}) \geq s_{\min}(\mathbf{A}) \geq \sqrt{n} - 2^{3.5} K^2 (\sqrt{\ln(9)d} + \sqrt{x})$$

with probability at least  $1 - e^{-x}$  for  $x > 0$ . Observing that  $\|\Sigma^{\dagger \frac{1}{2}}\| = s_{\min}^+(\Sigma)^{-1/2}$ , this implies that

$$s_{\min}(\mathbf{D}) \geq \sqrt{s_{\min}^+(\Sigma)} \left( \sqrt{n} - 2^{3.5} K^2 (\sqrt{\ln(9)d} + \sqrt{x}) \right),$$

while dividing through by  $\sqrt{n}$ , taking the non-negative part of the r.h.s. and squaring gives us

$$\lambda_{\min}(\widehat{\Sigma}) \geq \lambda_{\min}^+(\Sigma) \left( 1 - 2^{3.5} K^2 \left( \sqrt{\ln(9) \frac{d}{n}} + \sqrt{\frac{x}{n}} \right) \right)_+^2.$$

**Case  $n < d$ .** In this case we essentially study concentration of a smallest singular value of a Gram matrix  $\widehat{\mathbf{G}} = \frac{1}{d} \mathbf{D} \mathbf{D}^\top$ . For the case  $n < d$ , Theorem 4 would give us a vacuous estimate, and therefore we rely on Theorem 5 which requires additional assumption that columns of  $\mathbf{D}^\top$  lie on a (elliptic)

sphere of radius  $\sqrt{d}$ . In particular, similarly as before, applying Theorem 5 to the matrix  $\Sigma^{\dagger\frac{1}{2}} \mathbf{D}^\top$  with isotropic columns  $\Sigma^{\dagger\frac{1}{2}} \mathbf{X}_i$  satisfying  $\|\Sigma^{\dagger\frac{1}{2}} \mathbf{X}_i\| = \sqrt{d}$  a.s. for all  $i \in [n]$ , we get

$$\|\Sigma^{\dagger\frac{1}{2}}\| s_{\min}(\mathbf{D}^\top) \geq s_{\min}(\Sigma^{\dagger\frac{1}{2}} \mathbf{D}^\top) \geq \sqrt{d} - 2^{3.5} K^2 (\sqrt{\ln(9)n} + 6\sqrt{x})$$

with probability at least  $1 - e^{-x}$  for  $x > 0$ . Again, this gives us

$$s_{\min}(\mathbf{D}) \geq \sqrt{s_{\min}^+(\Sigma)} \left( \sqrt{d} - 2^{3.5} K^2 (\sqrt{\ln(9)n} + 6\sqrt{x}) \right),$$

while dividing through by  $\sqrt{d}$ , taking the non-negative part of the r.h.s. and squaring gives us

$$\lambda_{\min}(\widehat{\mathbf{G}}) \geq \lambda_{\min}^+(\Sigma) \left( 1 - 2^{3.5} K^2 \left( \sqrt{\ln(9) \frac{n}{d}} + 6\sqrt{\frac{x}{d}} \right) \right)_+^2.$$

Now we relate  $\lambda_{\min}(\widehat{\mathbf{G}})$  to the smallest non-zero eigenvalue of  $\widehat{\Sigma}$  (see also [Bai and Yin, 1993, Remark 1]). The smallest eigenvalue of  $d\widehat{\mathbf{G}}$  corresponds to  $d - n + 1$ -th smallest eigenvalue of  $n\widehat{\Sigma}$ , that is  $d\lambda_{\min}(\widehat{\mathbf{G}}) = n\lambda_{\min}^+(\widehat{\Sigma})$ . That said, multiplying the previous inequality through by  $d/n$  and rearranging, we get

$$\lambda_{\min}^+(\widehat{\Sigma}) \geq \lambda_{\min}^+(\Sigma) \left( \sqrt{\frac{d}{n}} - 2^{3.5} K^2 \left( \sqrt{\ln(9)} + 6\sqrt{\frac{x}{n}} \right) \right)_+^2$$

The proof is now complete.  $\square$

## E Minimum eigenvalue and condition number

Previous works considered the link between the condition number of the features and the DD behavior [Rangamani et al., 2020]. In this work, the analysis focuses more particularly on the minimum eigenvalue. In the following small experiments, we empirically show that in the experiments shown in the main paper, the condition number is driven by the minimum eigenvalue, and that the maximum eigenvalue stays close to a constant order when we increase the size of the features. In Figure 4, we use the same setting as in the MNIST experiment in Figure 2 is the main paper. We observe that the behavior of the condition number follows the minimum eigenvalue, while the maximum eigenvalue stays between 10 and 100 as we increase the width of the networks.

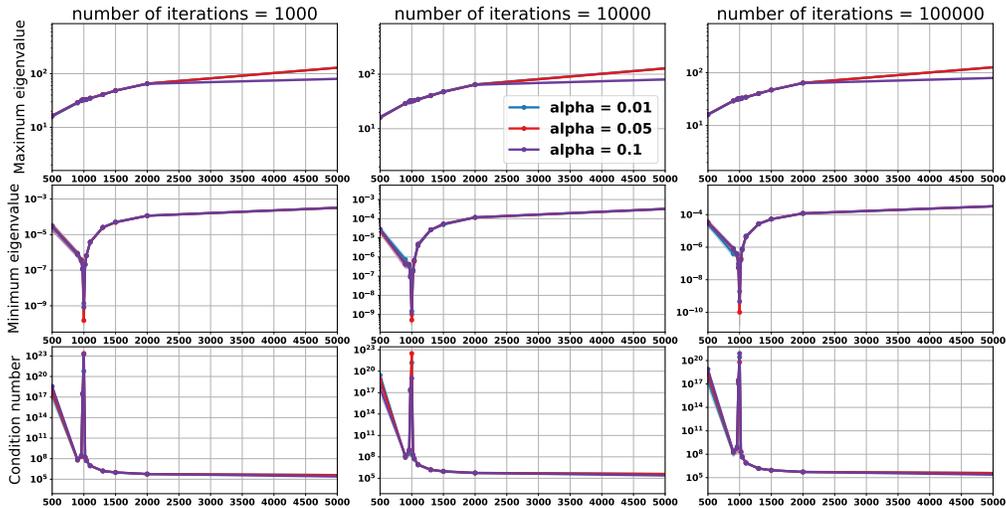


Figure 4: Maximum and minimum eigenvalues and condition numbers of the features of one hidden layer networks of variable width: MNIST - 1000 samples for training, networks trained with gradient descent and different step sizes.

## F More on the effect of depth

In section 6, we suggested that the ill-conditioning of the intermediary features of a neural network is not only due to the size of the network, but also to the weights distribution across the layers. More particularly, we suggest here that the optimization difficulty we observe for deep neural networks is linked among other factors to the minimum eigenvalue of the activations of the penultimate layer.

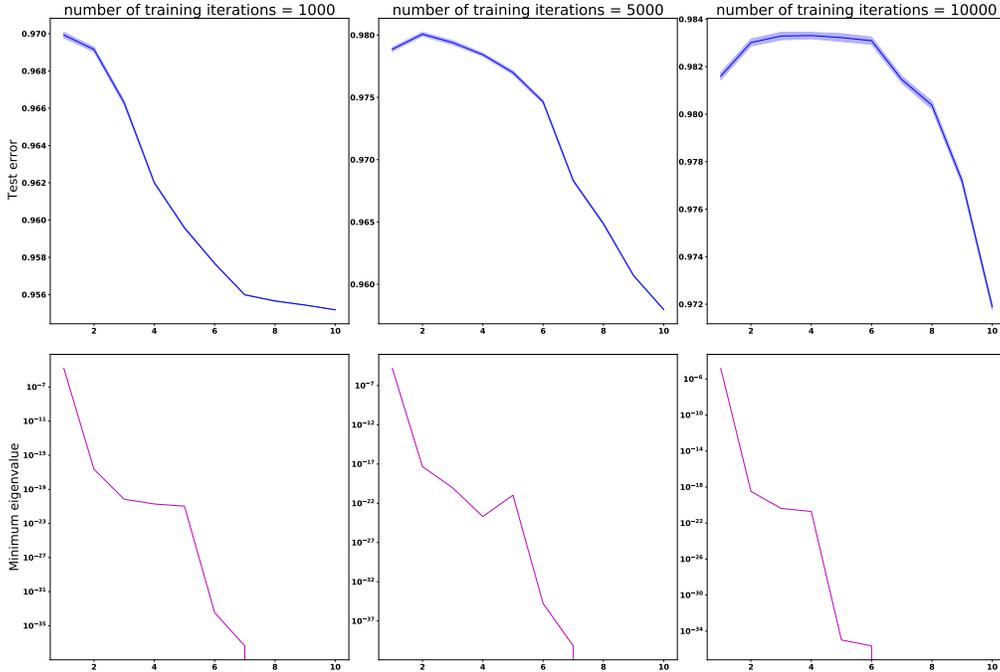


Figure 5: Mean test error and minimum eigenvalue for networks of fixed width = 500 and varying depth: MNIST - 1000 samples for training, 10000 samples for test, networks trained with gradient descent and step size 0.01.

To support our hypothesis, we run an experiment where we train networks of a fixed width (equal to 500) and depth varying from 2 to 10. We track the test error at various stages of training and the minimum eigenvalue of the features of the last layer. In Figure 5, we can observe that as expected, the deeper the network, the harder it is to train them. This is reflected in the increasing test error. For the deepest network, simple gradient descent fails to obtain a reasonable performance even after 10000 iterations. Moreover, we observe that the deeper the network, the smaller is the minimum eigenvalue, and the most ill-conditioned settings get even worse with training.

To further this analysis, we also compare networks with 3 hidden layers where we increase the width in all the layers and in the penultimate layer only, creating bottleneck in the earlier layers. This experiment complements Figure 3. In Figure 6, we observe that the bottleneck results in a more important drop in the eigenvalue around the width 1000 (width of the last layer in this case). Moreover, the minimum eigenvalue stays smaller than the other considered architectures when we increase the depth. This is reflected in a higher test error, confirming once more the effect of the conditioning of the last layer features on the final performance of the network when trained with gradient descent.

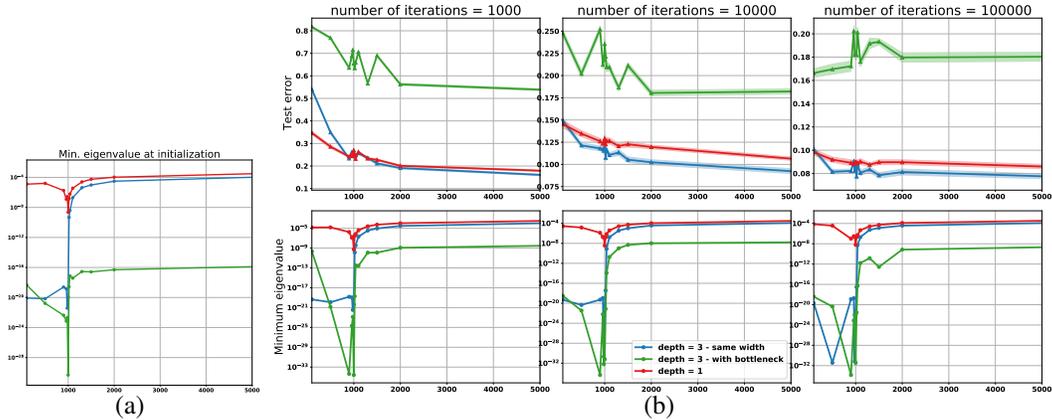


Figure 6: Training networks of increasing width with 1 and 3 hidden layers on MNIST - For the version with bottleneck, only the size of the last hidden layer is increased, while the other layers are composed of 10 neurons: (a) Minimum positive eigenvalue of the intermediary features at initialization - (b) Test error and corresponding minimum eigenvalue of the intermediary features at different iterations

## G Additional Empirical Evaluation

### G.1 Experimental settings - More details

In our experiments, we considered two datasets: MNIST and FashionMNIST. Both datasets have an input dimension of 784 and a training set of  $6 \cdot 10^4$  samples. As our theory predicts that the drop in the minimum eigenvalue and the performance of the models happens when the feature size reaches the size of the training set, and in order to keep our model tractable, we use subsets of size 1000 of the training sets. These subsets are randomly chosen and kept the same when the size of the model increases. All the models are trained with plain gradient descent, with a fixed step size. We use a step size of 0.01 unless stated otherwise. All the weights of the networks are initialized from a truncated normal distribution with a scaled variance. Finally, for the MNIST experiment in Figure 2, the mean and standard errors are estimated from runs with different seeds. For the other experiments, the mean and standard errors of the test error are estimated by splitting the test set into 10 subsets.

### G.2 More on the effect of architectural choices

In section 6, we suggested that for neural networks the quantity of interest might also be  $\hat{\lambda}_{\min}^+$  for intermediary features, which is affected by size of the model but also by the distribution of the weights and architectural choices. Section F shows some experiments that validate this hypothesis. To further our analysis, we question here the impact of skip connections on the eigenvalue of features at initialization. The difficulty that depth cause for the optimization of neural networks led to our reliance on skip connections among other tricks [De and Smith, 2020]. Here, we hypothesize that skip connections make the optimization of deep networks easier thanks to a better conditioning of the feature, through a less severe drop in the minimum eigenvalue around the interpolation threshold. Figure 7 shows that for a deep network with skip connection, the minimum eigenvalue of the penultimate layer activations behaves like this of a shallow neural network.

Eigenvalue of features at initialization with and without skip connection

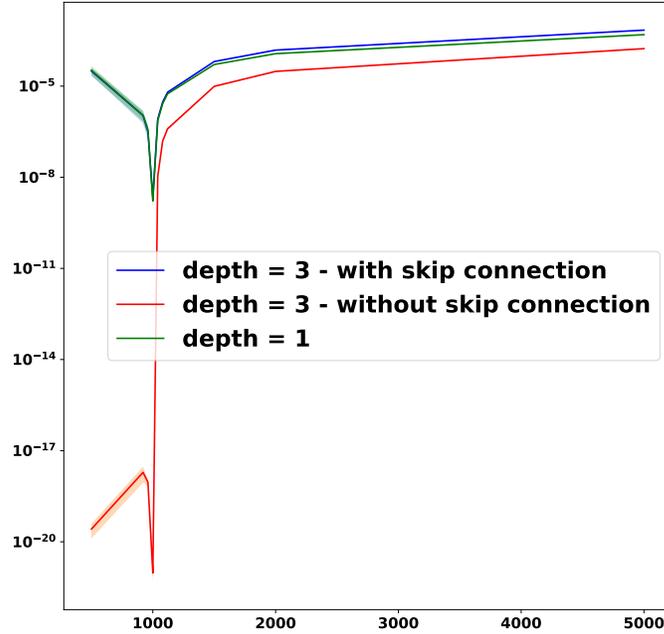


Figure 7: Mean minimum eigenvalue at initialization for networks of depths 1 and 3 and varying width. For the network of depth 3, we show two variants: with and without skip connection. The skip connection makes the deep network eigenvalue behave like the shallow network’s.

### G.3 On the position of the peak

The jump of the test error predicted by our theory and observed in our experiments appears when the width of the penultimate layer of the network reaches the size of the training set. The jump of the test loss observed in [Belkin et al. \[2019\]](#) appears when the size of the model reaches the size of the training set multiplied by the number of classes. It is then a natural question to ask whether these two positions coincide. Let us consider a network with a single hidden layer for simplicity. For an input of dimension  $d$ , a hidden layer of width  $h$  and a problem with  $k$  classes, the size of such a network is  $(d + 1)h + (h + 1)k$ . For a training set of size  $n$ , [Belkin et al. \[2019\]](#) observes a peak in the test loss when  $(d + 1)h + (h + 1)k = nk$ , which corresponds to a model with a hidden layer of width  $h = \frac{nk - k}{d + k + 1}$ . When training with 1000 samples from MNIST (as in our experiments), this width is then  $h \approx 12$ . Recall that the peak highlighted here appears at  $h = 1000$ , which is a much bigger model. This observation reinforces the observation that neural networks show not only a double but a triple [[d’Ascoli et al., 2020](#)] (or a multiple?) descent behavior. A thorough verification is however beyond the scope of this work, and is left to future research.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [\[Yes\]](#)
  - (b) Did you describe the limitations of your work? [\[Yes\]](#) These are summarised in Section 7
  - (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#) We believe that presented research should be categorized as basic research and we are not targeting any specific application area. Theorems may inspire new algorithms and theoretical investigation. The algorithms presented here can be used for many different applications and a particular use may have both positive or negative impacts. We are not aware of any immediate short term negative implications of this research and we believe that a broader impact statement is not required for this paper.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) Assumptions are stated in all statements of theorems.
  - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) The main paper presents proof sketch while all the proofs with details are deferred to the supplementary material.
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) Instructions for experiments are presented in Section 6
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#)
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#)
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) Synthetic experiments presented in the introduction can be run on any contemporary laptop, resources are also discussed in Section 6.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#)
  - (b) Did you mention the license of the assets? [\[N/A\]](#)
  - (c) Did you include any new assets either in the supplemental material or as a URL? [\[N/A\]](#)
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [\[N/A\]](#)
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#)
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)