



HAL
open science

Phase diagram of Stochastic Gradient Descent in high-dimensional two-layer neural networks

Rodrigo Veiga, Ludovic Stephan, Bruno Loureiro, Florent Krzakala, Lenka
Zdeborová

► **To cite this version:**

Rodrigo Veiga, Ludovic Stephan, Bruno Loureiro, Florent Krzakala, Lenka Zdeborová. Phase diagram of Stochastic Gradient Descent in high-dimensional two-layer neural networks. 2022. hal-04026190

HAL Id: hal-04026190

<https://hal.science/hal-04026190>

Preprint submitted on 13 Mar 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Phase diagram of Stochastic Gradient Descent in high-dimensional two-layer neural networks

Rodrigo Veiga^{1,3}, Ludovic Stephan¹, Bruno Loureiro¹, Florent Krzakala¹, and Lenka Zdeborová²

¹Ecole Polytechnique Fédérale de Lausanne (EPFL). Information, Learning and Physics (IdePHICS) lab.
CH-1015 Lausanne, Switzerland.

²Ecole Polytechnique Fédérale de Lausanne (EPFL). Statistical Physics of Computation (SPOC) lab.
CH-1015 Lausanne, Switzerland.

³Universidade de São Paulo. Instituto de Física. São Paulo, SP, Brazil.

Abstract

Despite the non-convex optimization landscape, over-parametrized shallow networks are able to achieve global convergence under gradient descent. The picture can be radically different for narrow networks, which tend to get stuck in badly-generalizing local minima. Here we investigate the cross-over between these two regimes in the high-dimensional setting, and in particular investigate the connection between the so-called mean-field/hydrodynamic regime and the seminal approach of Saad & Solla. Focusing on the case of Gaussian data, we study the interplay between the learning rate, the time scale, and the number of hidden units in the high-dimensional dynamics of stochastic gradient descent (SGD). Our work builds on a deterministic description of SGD in high-dimensions from statistical physics, which we extend and for which we provide rigorous convergence rates.

Contents

1	Introduction	3
2	Setting	5
3	Main results	7
4	Discussion, special cases, and simulations	11
4.1	Saad & Solla scaling $\kappa = \delta = 0$	11
4.2	Perfect learning for $\kappa = 0$	12
4.3	Bad learning for $\kappa = 0$	13
4.4	Large hidden layer: $\kappa > 0$	14
5	Conclusion	14
	Appendix	16
A	Deterministic scaling limit of stochastic processes	16
A.1	Preliminaries: bounding the q_{jj}	17
A.2	Assumption A.1.1	18
A.3	Assumption A.1.2	19
A.4	$\sqrt{\cdot}$ -Lipschitz property	19
B	A lemma on ODE perturbation	19
C	Expectations over the local fields	21
C.1	Population risk	21
C.2	ODE contributions	22
C.3	From gradient flow to local fields	24
D	Initial conditions and symmetric teacher	24
	References	26

1 Introduction

Descent-based algorithms such as stochastic gradient descent (SGD) and its variants are the workhorse of modern machine learning. They are simple to implement, efficient to run and most importantly: they work well in practice. A detailed understanding of the performance of SGD is a major topic in machine learning. Quite recently, significant progress was achieved in the context of learning in shallow neural networks. In a series of works, it was shown that the optimisation of wide two-layer neural networks can be mapped to a convex problem in the space of probability distributions over the weights [1, 2, 3, 4]. This remarkable result implies global convergence of two-layer networks towards perfect learning provided that the number of hidden neurons is large, the learning rate is sufficiently small and enough data is at disposition. This line of work is commonly referred to as the mean-field or the hydrodynamic limit of neural networks. Mathematically, these works showed that one could describe the entire dynamics using a partial differential equation (PDE) in d dimensions.

In a different, and older, line of work one-pass SGD for two-layer neural networks with a *finite number* p of hidden units, synthetic Gaussian input data and teacher-generated labels has been widely studied starting with the seminal work of [5]. These works consider the limit of high-dimensional data and show, in particular, that the stochastic process driven by gradient updates converge to a set of p^2 deterministic ordinary differential equations (ODEs) as the input dimension $d \rightarrow \infty$ and the learning rate is proportional to $1/d$. The validity of these ODEs in this limit was proven by [6]. However, the picture drawn from the analysis of these ODEs is slightly different from the mean-field/hydrodynamic picture: in this case SGD can get stuck for long time in minima associated to no specialization of the hidden units to the teacher hidden units, and even when it converges to specializing minima, it fails to perfectly learn (i.e. to achieve zero population risk). In fact, in this analysis, the interplay between the limit of the learning rate going to zero and $d \rightarrow \infty$ appeared to be fundamental.

One should naturally wonder about the link between these two sets of works with, on the one hand a d -dimensional PDE (with large p), and on the other a p^2 -dimensional ODE (with large d). In this work we aim to build a bridge between these two approaches for studying one-pass SGD.

Our starting point is the framework from [5], which we build upon and expand to a much broader range of choices of learning rate, time scales, and hidden layer width. This allows us to provide a sharp characterisation of the performance of SGD for two-layer neural networks in high-dimensions. We show it depends on the precise way in which the limit is taken, and in particular on how the quantity of data, the hidden layer width, and the learning rate scale as $d \rightarrow \infty$. For different choices of scaling, we can observe scenarios such as perfect learning, imperfect learning with an unavoidable error, or even no learning at all.

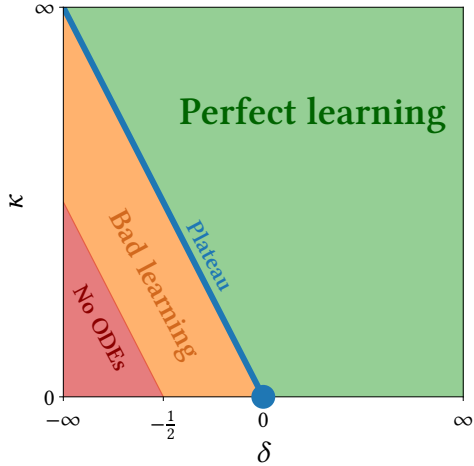
As a consequence of our analysis, we provide a phase diagram (see Figure 1a) describing the possible scenarios arising in the high-dimensional setting. Our **main contributions** are as follow:

C1 We rigorously show that the dynamics of SGD can be captured by a set of deterministic ODEs, considerably extending the proof of [6] to accommodate for general time scalings defined by an arbitrary learning rate, and a general range of hidden layer width. We provide much finer non-asymptotic guarantees which are crucial for our subsequent analysis.

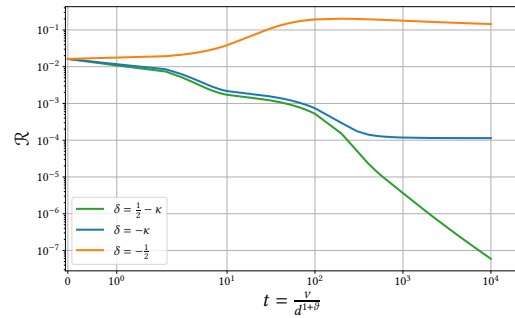
C2 From the analysis of the ODEs, we derive a phase diagram of SGD for two-layer neural networks in the high-dimensional input layer limit $d \rightarrow \infty$. In particular, scaling both the learning rate γ and hidden layer width p with the input dimension d as

$$\gamma \propto d^{-\delta}, \tag{1a}$$

$$p \propto d^\kappa, \tag{1b}$$



(a) The phase diagram of SGD learning regimes for two-layer neural networks in the high-dimensional input layer limit $d \rightarrow \infty$. Eqs. (1) define proper time scalings for each of the regions. Perfect learning region: $\kappa + \delta > 0$. Plateau line: $\kappa + \delta = 0$. Bad learning region: $-1/2 < \kappa + \delta < 0$. No ODEs region: $\kappa + \delta < -1/2$.



(b) A solution of the ODEs in all regions of Figure 1a, with matching colors. Parameters $\kappa = 0.301$, $p = 8$, $k = 4$, $\rho_{rs} = \delta_{rs}$. Noise: $\Delta = 10^{-3}$. Activation function: $\sigma(x) = \text{erf}(x/\sqrt{2})$. Data distribution: $\mathbb{P}(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbb{1})$. The time scaling is not uniform through the phase diagram: $\vartheta = \kappa + \delta$ on green and blue regimes and $\vartheta = 2(\kappa + \delta)$ on the orange region. The green curve decays as a power law to zero excess error.

Figure 1: Phase diagram (left) and typical behavior of the ODE in each regions (right).

we identify four different learning regimes which are summarized in Figure 1a:

- Perfect learning (green region, $\kappa > -\delta$): we show that perfect learning (zero population risk) can be asymptotically achieved with $n \sim d^{1+\kappa+\delta}$ samples even for tasks with additive noise.
- Plateau (blue line $\kappa = -\delta$): learning reaches a plateau related to the noise strength. The point $\kappa = \delta = 0$ goes back to the classical work of [5].
- Bad learning (orange region $-1/2 < \kappa + \delta < 0$): here the noise dominates the learning process.
- No ODEs (red region $\kappa + \delta < -1/2$): the stochastic process associated to SGD is not guaranteed to converge to a set of deterministic ODEs. This region is thus outside the scope of our analysis.

To better illustrate this phase diagram we present in Figure 1b a solution of the ODEs in all three regimes.

Relation to previous work – Deterministic dynamical descriptions of one-pass stochastic gradient descent in high-dimensions have a long tradition in the statistical physics community, starting with single- and two-layer neural networks with few hidden units [7, 8, 9, 10, 11]. The seminal work by [5] overcame previous limitations by constructing a set of deterministic ODEs for two-layer networks with any finite number of hidden units, paving the way for a series of important contributions [12, 13, 14, 6]. This line of work corresponds to the $\kappa = \delta = 0$ case of Figure 1a. One of our goal is to generalize this picture beyond fixed hidden layer size and learning rate.

A more recent line of work investigating the dynamics of SGD is the so-called *mean-field limit* [1, 15, 2, 3, 4], which connects the SGD dynamics of large-width two-layer neural networks to a diffusion equation

in the hidden layer weight density. In particular, [15] provide non-asymptotic convergence bounds for sufficiently small learning rates, corresponding to the green region of Figure 1a (with $p \rightarrow \infty$). The mean-field approach computes the empirical distribution (in \mathbb{R}^d) of the hidden layer weights, while we focus on the macroscopic overlaps between the teacher and student weights.

Reproducibility A code is provided at https://github.com/rodsveiga/phdiag_sgd.

2 Setting

Consider a supervised learning regression task. The data set is composed of n pairs $(\mathbf{x}^v, y^v)_{v \in [n]} \in \mathbb{R}^{d+1}$ identically and independently sampled from $\mathbb{P}(\mathbf{x}, y)$. The probability $\mathbb{P}(\mathbf{x})$ is assumed to be known and $\mathbb{P}(y|\mathbf{x})$ is modelled by a two layer neural network called the *teacher*. Given a feature vector $\mathbf{x}^v \in \mathbb{R}^d$, the respective label $y^v \in \mathbb{R}$ is defined as the output of a network with k hidden units, fixed weights $\mathbf{W}^* \in \mathbb{R}^{k \times d}$ and an activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$:

$$y^v = f(\mathbf{x}^v, \mathbf{W}^*) + \sqrt{\Delta} \zeta^v, \quad (2)$$

where

$$f(\mathbf{x}^v, \mathbf{W}^*) = \frac{1}{k} \sum_{r=1}^k \sigma \left(\frac{\mathbf{w}_r^{*\top} \mathbf{x}^v}{\sqrt{d}} \right) = \frac{1}{k} \sum_{r=1}^k \sigma(\lambda_r^{*v}), \quad (3)$$

with $\mathbf{w}_r^* \equiv [\mathbf{W}^*]_r \in \mathbb{R}^d$ as the r -th row of the matrix \mathbf{W}^* and $\lambda_r^{*v} \equiv \mathbf{w}_r^{*\top} \mathbf{x}^v / \sqrt{d} \in \mathbb{R}$ as the r -th component of the teacher *local field* vector $\boldsymbol{\lambda}^{*v} \in \mathbb{R}^k$. The parameter $\Delta \geq 0$ controls the strength of additive label noise: $\zeta^v \sim \mathbb{P}(\zeta^v)$ such that $\mathbb{E}_{\zeta \sim \mathbb{P}(\zeta)}[\zeta] = 0$ and $\mathbb{E}_{\zeta \sim \mathbb{P}(\zeta)}[\zeta^2] = 1$.

Given a new sample $\mathbf{x} \sim \mathbb{P}(\mathbf{x})$ outside the training data, the goal is to obtain an estimation $\hat{f}(\mathbf{x})$ for the respective label y . The error is quantified by a loss function $\mathcal{L}(y, \hat{f}(\mathbf{x}, \boldsymbol{\Theta}))$, where $\boldsymbol{\Theta}$ is an arbitrary set of parameters to be learned from data.

In this manuscript we are interested in the problem of estimating \mathbf{W}^* with another two-layer neural network with the same activation function, which we will refer to as the *student*. The student network has p hidden units and a matrix of weights $\mathbf{W} \in \mathbb{R}^{p \times d}$ to be *learned* from the data. Given a feature vector $\mathbf{x} \sim \mathbb{P}(\mathbf{x})$ the student prediction for the respective label is given as

$$\hat{f}(\mathbf{x}, \mathbf{W}) = \frac{1}{p} \sum_{j=1}^p \sigma \left(\frac{\mathbf{w}_j^\top \mathbf{x}}{\sqrt{d}} \right) = \frac{1}{p} \sum_{j=1}^p \sigma(\lambda_j^v), \quad (4)$$

where $\mathbf{w}_j \equiv [\mathbf{W}]_j \in \mathbb{R}^d$ is the j -th row of the matrix \mathbf{W} and $\lambda_j \equiv \mathbf{w}_j^\top \mathbf{x} / \sqrt{d} \in \mathbb{R}$ is defined as j -th component of the student *local field* vector $\boldsymbol{\lambda} \in \mathbb{R}^p$.

One-pass gradient descent – Typically, one minimizes the *empirical risk* over the full data set. Instead, learning with *one-pass* gradient descent minimizes directly the *population risk*:

$$\mathcal{R}(\mathbf{W}, \mathbf{W}^*) \equiv \mathbb{E}_{\mathbf{x}, y \sim \mathbb{P}(\mathbf{x}, y)} \left[\mathcal{L} \left(f(\mathbf{x}, \mathbf{W}^*), \hat{f}(\mathbf{x}, \mathbf{W}) \right) \right]. \quad (5)$$

Given a *single* sample (\mathbf{x}^v, y^v) the weights are updated sequentially by the gradient descent rule:

$$\mathbf{w}_j^{v+1} = \mathbf{w}_j^v - \gamma \nabla_{\mathbf{w}_j} \mathcal{L} \left(y^v, \hat{f}(\mathbf{x}^v, \mathbf{W}) \right), \quad (6)$$

with $\nu \in [n]$ and $j \in [p]$. The parameter $\gamma > 0$ is the learning rate. Despite being a simplification with respect to batch learning, one-pass gradient descent is an amenable surrogate for the theoretical analysis of non-convex optimization, since at each step the gradient is computed with a fresh data sample, which is equivalent to performing SGD directly on the population risk.

In particular, in this manuscript we assume realizability $p \geq k$, and focus our analysis on the square loss $\mathcal{L}(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$, leading to

$$\mathbf{w}_j^{\nu+1} = \mathbf{w}_j^\nu + \frac{\gamma}{p\sqrt{d}} \sigma'(\lambda_j^\nu) \mathcal{E}^\nu \mathbf{x}^\nu, \quad (7)$$

where

$$\mathcal{E}^\nu \equiv \frac{1}{k} \sum_{r=1}^k \sigma(\lambda_r^{*\nu}) - \frac{1}{p} \sum_{l=1}^p \sigma(\lambda_l^\nu) + \sqrt{\Delta} \zeta^\nu. \quad (8)$$

with population risk given by

$$\mathcal{R}(\mathbf{W}, \mathbf{W}^*) = \frac{1}{2} \mathbb{E}_{\mathbf{x}, y \sim \mathbb{P}(\mathbf{x}, y)} \left[\left(\hat{f}(\mathbf{x}, \mathbf{W}) - f(\mathbf{x}, \mathbf{W}^*) \right)^2 \right]. \quad (9)$$

Therefore, from the above expression we can see that to monitor the population risk along the learning dynamics it is sufficient to track the joint distribution of the local fields $(\boldsymbol{\lambda}, \boldsymbol{\lambda}^*)$. For Gaussian data $\mathbb{P}(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbf{1})$, one can replace the expectation $\mathbb{E}_{\mathbf{x}, y \sim \mathbb{P}(\mathbf{x}, y)}[\cdot]$ by $\mathbb{E}_{\boldsymbol{\lambda}, \boldsymbol{\lambda}^* \sim \mathcal{N}(\boldsymbol{\lambda}, \boldsymbol{\lambda}^*|\mathbf{0}, \Omega)}[\cdot]$ and fully describe the dynamics through the following sufficient statistics, known in the statistical physics literature as *macroscopic variables*:

$$\mathbf{Q}^\nu \equiv \mathbb{E}_{\mathbf{x}, y \sim \mathbb{P}(\mathbf{x}, y)} [\boldsymbol{\lambda}^\nu \boldsymbol{\lambda}^{\nu\top}] = \frac{1}{d} \mathbf{W}^{\nu\top} \mathbf{W}^\nu, \quad (10a)$$

$$\mathbf{M}^\nu \equiv \mathbb{E}_{\mathbf{x}, y \sim \mathbb{P}(\mathbf{x}, y)} [\boldsymbol{\lambda}^\nu \boldsymbol{\lambda}^{*\nu\top}] = \frac{1}{d} \mathbf{W}^{\nu\top} \mathbf{W}^*, \quad (10b)$$

$$\mathbf{P} \equiv \mathbb{E}_{\mathbf{x}, y \sim \mathbb{P}(\mathbf{x}, y)} [\boldsymbol{\lambda}^{*\nu} \boldsymbol{\lambda}^{*\nu\top}] = \frac{1}{d} \mathbf{W}^{*\top} \mathbf{W}^*. \quad (10c)$$

with matrix elements, called *order parameters* in the statistical physics literature, denoted by $q_{jl}^\nu \equiv [\mathbf{Q}^\nu]_{jl}$, $m_{jr}^\nu \equiv [\mathbf{M}^\nu]_{jr}$ and $\rho_{rs} \equiv [\mathbf{P}]_{rs}$. The *macroscopic state* of the system at the learning step ν is given by the *overlap matrix* $\boldsymbol{\Omega}^\nu \in \mathbb{R}^{(p+k) \times (p+k)}$:

$$\boldsymbol{\Omega}^\nu = \begin{bmatrix} \mathbf{Q}^\nu & \mathbf{M}^\nu \\ \mathbf{M}^{\nu\top} & \mathbf{P} \end{bmatrix}, \quad (11)$$

and the population risk is completely determined by the macroscopic state:

$$\mathcal{R}(\boldsymbol{\Omega}) = \frac{1}{2} \mathbb{E}_{\boldsymbol{\lambda}, \boldsymbol{\lambda}^* \sim \mathcal{N}(\boldsymbol{\lambda}, \boldsymbol{\lambda}^*|\mathbf{0}, \Omega)} \mathbb{E}_{\zeta \sim \mathbb{P}(\zeta)} \left[\left(\hat{f}(\boldsymbol{\lambda}) - f(\boldsymbol{\lambda}^*) \right)^2 \right]. \quad (12)$$

The training dynamics (6) defines a discrete-time stochastic process for the evolution of the overlap matrix

$$\left\{ \boldsymbol{\Omega}^\nu \in \mathbb{R}^{(p+k) \times (p+k)}, \nu \in [n] \right\}, \quad (13)$$

with \mathbf{P} fixed and \mathbf{Q}^ν and \mathbf{M}^ν updated as:

$$q_{jl}^{\nu+1} - q_{jl}^\nu = \frac{\gamma}{pd} \underbrace{\left(\mathcal{E}_j^\nu \lambda_l^\nu + \mathcal{E}_l^\nu \lambda_j^\nu \right)}_{\text{learning}} + \frac{\gamma^2 \|\mathbf{x}\|^2}{p^2 d^2} \underbrace{\mathcal{E}_j^\nu \mathcal{E}_l^\nu}_{\text{variance}}, \quad (14a)$$

$$m_{jr}^{v+1} - m_{jr}^v = \frac{\gamma}{pd} \underbrace{\mathcal{E}_j^v \lambda_r^{*v}}_{\text{learning}}, \quad (14b)$$

with $v \in [n]$, $j, l \in [p]$, $r \in [k]$ and $\mathcal{E}_j^v \equiv \sigma'(\lambda_j^v) \mathcal{E}^v$. In what follows, we will make the concentration assumption $\|\mathbf{x}\|^2 = d$; this will be justified in the proof of Theorem 3.1.

We emphasize in (14) the specific role played by each term in the right hand-side. The "learning" terms are the fundamental ones, that actually drive the learning of the teacher by the student. We show in Appendix C.3 that these "learning" terms are identical to those obtained in the gradient flow approximation of SGD, whose performance is the topic of many works [1, 2, 3, 4]. Those are *precisely* the terms that draw the population risk towards zero. However, in our setting there is an additional variance term (so that this flow approximation is incomplete) that corresponds to the fluctuations of $\mathcal{L}(\mathbf{x}, \mathbf{W}, \mathbf{W}^*)$ around its expected value $\mathcal{R}(\mathbf{W}, \mathbf{W}^*)$. In particular, this is where the effects of the noise ζ can be felt. These terms were sometimes denoted as (I_2) and (I_4) in [16]. We shall see that the additional "variance" term is the one responsible for the plateau in the critical (blue) region of Figure 1a, while its contribution vanishes in the perfect learning (green) region.

Additionally, albeit our work particularizes to Gaussian input data, we believe our conclusion, and the phase diagram discussed in Figure 1a, to hold beyond this restricted case. Indeed, while the Gaussian assumption is crucial to reach a particular set of ODEs and their analytic expression, the approach can be applied to more complex data distribution, as long as one can track the sufficient statistics required to have a closed set of equations. For instance, [17] obtained very similar equations for an arbitrary mixture of Gaussians – that would obey the same scaling analysis as ours – while [18, 19, 20] proved that many complex distributions behave as Gaussians in high-dimensional setting, including, e.g. realistic GAN-generated data. We thus expect our conclusions to be robust in this respect.

3 Main results

Although $t_0 = v/d$ seems to be the most natural time scaling in the high-dimensional limit $d \rightarrow \infty$, if γ and p are allowed to vary with d the right-hand side (RHS) of Eqs. (14) can diverge and render the ODE approximation obsolete. Instead, for a given time scaling δt , we can rewrite Eqs. (14) as

$$\frac{q_{jl}^{v+1} - q_{jl}^v}{\delta t} = \frac{\gamma}{pd \delta t} \left(\mathcal{E}_j^v \lambda_l^v + \mathcal{E}_l^v \lambda_j^v \right) + \frac{\gamma^2}{p^2 d \delta t} \mathcal{E}_j^v \mathcal{E}_l^v, \quad (15a)$$

$$\frac{m_{jr}^{v+1} - m_{jr}^v}{\delta t} = \frac{\gamma}{pd \delta t} \mathcal{E}_j^v \lambda_r^{*v}. \quad (15b)$$

In Theorem 3.1 we prove that as $d \rightarrow \infty$, Ω^v converges to the solution of the ODE:

$$\frac{d}{dt} \bar{\Omega}(t) = \psi(\bar{\Omega}(t)), \quad (16)$$

where $\psi : \mathbb{R}^{(p+k) \times (p+k)} \rightarrow \mathbb{R}^{(p+k) \times (p+k)}$ is the expected value of the RHS of Eqs. (15), provided that this solution stays bounded. This enhances the result of [6] by providing convergence rates to the ODEs encompassing all scalings adopted hereafter:

Theorem 3.1 (Deterministic scaling limit of stochastic processes). *Let $\tau \in \mathbb{R}$ be the continuous time horizon and $\delta t = \delta t(d)$ be a time scaling factor such that the following assumptions hold:*

1. the time scaling δt satisfies for some constant c ,

$$\delta t \geq c \max\left(\frac{\gamma}{pd}, \frac{\gamma^2}{p^2d}\right) \quad (17)$$

2. the activation function σ is L -Lipschitz,

3. the function $\psi : \mathbb{R}^{(p+k) \times (p+k)} \rightarrow \mathbb{R}^{(p+k) \times (p+k)}$ is L' -Lipschitz.

Then, there exists a constant $C > 0$ (depending on c, L, L') such that for any $0 \leq \nu \leq \lfloor \tau/\delta t \rfloor$, the following inequality holds:

$$\mathbb{E} \|\Omega^\nu - \bar{\Omega}(\nu\delta t)\|_\infty \leq e^{C\nu} \log(p) \sqrt{\delta t}. \quad (18)$$

Our proof is based on techniques introduced in [21] (namely, their Lemma 2) which studies a different problem with related proof techniques. The proof involves decomposing $\Omega^{\nu+1}$ as

$$\Omega^{\nu+1} = \Omega^\nu + \delta t \psi(\Omega^\nu) + (\Omega^{\nu+1} - \Omega^\nu - \delta t \psi(\Omega^\nu)), \quad (19)$$

where the two first terms can be considered as a deterministic discrete process, and the last term is a martingale increment. The main challenge lies in showing that the martingale contribution stays bounded throughout the considered time period.

Although the method is similar to [6], there are a number of differences between the two approaches. First, our proof fixes a number of holes in [6], in particular bounding q_{jj}^ν by a sufficiently slowly diverging function of ν . Additionally, the techniques used in this paper yield a dependency in p that is nearly negligible, while the previous methods imply bounds that are much too coarse for our needs.

The function ψ can be computed explicitly for various choices of σ , which allows to check Assumption 3 directly. We provide in Appendix C the necessary computations for $\sigma(x) = \text{erf}(x/\sqrt{2})$; those for the ReLU unit can be found in [22]. It can be checked that in the ReLU case, the function ψ is not Lipschitz around the matrices Ω satisfying

$$\Omega_{jl} = \sqrt{\Omega_{jj}\Omega_{ll}}$$

for some $j \neq l$. However, in every case we have a weaker square-root-Lipschitz property: there exists $C \in \mathbb{R}$ such that

$$\|\psi(\Omega) - \psi(\Omega')\| \leq C \|\sqrt{\Omega} - \sqrt{\Omega'}\|$$

for any Ω, Ω' . Since the square root function is Lipschitz whenever the eigenvalues of Ω are bounded away from zero (see e.g. [23]), Assumption 3 is implied by the condition

$$\Omega^\nu \geq \epsilon I_{p+k};$$

however, this assumption is much stronger, and becomes unrealistic in the specialization phase (as well as when $p \gg d$).

Theorem 3.1 allows us to safely navigate through Figure 1a by keeping track of convergence rates of the discrete process to a set ODEs. The interplay between learning rate and hidden layer width defines the time scaling δt and the trade-off between the linear contribution on \mathcal{E}_j and the quadratic one, playing a central role on whether the network achieves perfect learning or not. Specifically, consider the following learning rate and hidden layer width scaling with d :

$$\gamma = \frac{\gamma_0}{d^\delta}, \quad (20a)$$

$$p = p_0 d^\kappa, \quad (20b)$$

where $\gamma_0 \in \mathbb{R}^+$ and $p_0 \in \mathbb{N}$ are constants. The exponent $\delta \in \mathbb{R}$ can be either greater or smaller than zero, while $\kappa \in \mathbb{R}^+$. Replacing these scalings on Eqs. (14), we find:

$$q_{jl}^{v+1} - q_{jl}^v = \frac{1}{d^{1+\kappa+\delta}} \underbrace{\left(\mathcal{E}_j^v \lambda_l^v + \mathcal{E}_l^v \lambda_j^v \right)}_{\text{learning}} + \frac{1}{d^{1+2(\kappa+\delta)}} \underbrace{\mathcal{E}_j^v \mathcal{E}_l^v}_{\text{noise}}, \quad (21a)$$

$$m_{jr}^{v+1} - m_{jr}^v = \frac{1}{d^{1+\kappa+\delta}} \underbrace{\mathcal{E}_j^v \lambda_r^{*v}}_{\text{learning}}, \quad (21b)$$

where we have chosen $\gamma_0 = p_0$ without loss of generality.

Since the distribution of the label noise $\mathbb{P}(\zeta)$ is such that $\mathbb{E}_{\zeta \sim \mathbb{P}(\zeta)}[\zeta] = 0$, the linear contribution in \mathcal{E}_j is noiseless in the high-dimensional limit $d \rightarrow \infty$, and therefore we will refer to it as the *learning term*. The noise enters in the equations through the variance computed on the quadratic contribution $\mathcal{E}_j \mathcal{E}_l$, which we will refer to as the *noise term*; intuitively, it is a high-dimensional variance correction which hinders learning. In order to satisfy (17), we shall take

$$\delta t = \max \left(\frac{1}{d^{1+\kappa+\delta}}, \frac{1}{d^{1+2(\kappa+\delta)}} \right). \quad (22)$$

When $\kappa + \delta \neq 0$, this implies that either the learning term or the noise term scale like a negative power of d , and is negligible with respect to the other term. It is then easy to check that at a finite time horizon τ , the resulting ODEs behave as if the negligible term was not present. We refer to Theorem B.1 in the appendix for a quantitative proof of this phenomenon. Let us now describe the different regimes depicted in Figure 1a.

Blue line (plateau) – When γ and p are scaled such that $\kappa = -\delta$, Eqs. (21) converge to

$$\frac{dq_{jl}}{dt_0} = \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \left[\mathcal{E}_j \lambda_l + \mathcal{E}_l \lambda_j \right] + \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \mathbb{E}_{\zeta \sim \mathbb{P}(\zeta)} \left[\mathcal{E}_j \mathcal{E}_l \right], \quad (23a)$$

$$\frac{dm_{jr}}{dt_0} = \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \left[\mathcal{E}_j \lambda_r^* \right], \quad (23b)$$

with $\delta t_0 \equiv 1/d$. This regime is an extension of [5] for which $\kappa = \delta = 0$. The convergence rate to the ODEs scales with $d^{-1/2} \log(d)$, and the phenomenology we observe for $\kappa = \delta = 0$ is consistent with previous works studying the setting $\kappa = \delta = 0$; namely the existence of an asymptotic plateau proportional to the noise level. For instance, the asymptotic population risk \mathcal{R}_∞ is known to be proportional to $\gamma \Delta$ [6] when $\kappa = \delta = 0$ and the dynamics is driven by a rescaled version of Eqs. (23). Since the noise term does not vanish under this scaling, perfect learning to zero population risk is not possible. There is always an asymptotic plateau related to the noise level Δ , and the learning rate γ .

Green region (perfect learning) – If $\kappa > -\delta$ we can define the time scaling $\delta t_{\kappa+\delta} \equiv 1/d^{1+\kappa+\delta}$. By Theorem 3.1, Eqs. (21) converge to the following deterministic set of ODEs:

$$\frac{dq_{jl}}{dt_{\kappa+\delta}} = \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \left[\mathcal{E}_j \lambda_l + \mathcal{E}_l \lambda_j \right] + \mathcal{O} \left(\frac{\mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \mathbb{E}_{\zeta \sim \mathbb{P}(\zeta)} \left[\mathcal{E}_j \mathcal{E}_l \right]}{d^{\kappa+\delta}} \right), \quad (24a)$$

$$\frac{dm_{jr}}{dt_{\kappa+\delta}} = \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} [\mathcal{E}_j \lambda_r^*] , \quad (24b)$$

at a rate proportional to $d^{-(1+\kappa+\delta)/2} \log(d)$, where we have highlighted that the noise term vanishes with $d^{-(\kappa+\delta)}$. Hence, as long as $\kappa > -\delta$ the noise does not play any role on the dynamics. This setting could be understood by taking an effect learning rate $\gamma_{\text{eff}} \propto d^{-\kappa-\delta}$ on $\mathcal{R}_\infty \propto \gamma \Delta$, which leads to zero population risk, i.e. perfect learning, in the high dimensional limit $d \rightarrow \infty$. We validate this claim by a finite size analysis in the next section.

As discussed, the time scaling determines the number of data samples required to complete one learning step on the continuous scale. The bigger $\kappa+\delta$, the more attenuated the noise term, thus the closer to perfect learning. The trade-off is that the bigger $\kappa+\delta$, the larger the number of samples needed is, since $n = \tau d^{1+\kappa+\delta}$. Given a realizable learning task, one would thus rather choose the parameters to attain the perfect learning region, but being as close as possible to the plateau line for not increasing too much the needed number of samples. We remark that [15] provides an alternative deterministic approximation in this regime, with non-asymptotic bounds, whenever $p \gg 1$; this is the so-called mean-field approximation, with known convergence guarantees [2].

Orange region (bad learning) – We now step in the unusual situation where the learning rate grows faster with d than the hidden layer width: $\kappa < -\delta$. In this case, by (22) the noise term dominates over the dynamics. Defining the time scaling $\delta t_{2(\kappa+\delta)} \equiv 1/d^{1+2(\kappa+\delta)}$, we have

$$\frac{dq_{jt}}{\delta t_{2(\kappa+\delta)}} = \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \mathbb{E}_{\zeta \sim \mathbb{P}(\zeta)} [\mathcal{E}_j \mathcal{E}_l] + \mathcal{O} \left(\frac{\mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} [\mathcal{E}_j \lambda_l + \mathcal{E}_l \lambda_j]}{d^{-(\kappa+\delta)}} \right) , \quad (25a)$$

$$\frac{dm_{jr}}{\delta t_{2(\kappa+\delta)}} = \mathcal{O} \left(\frac{\mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} [\mathcal{E}_j \lambda_r^*]}{d^{-(\kappa+\delta)}} \right) . \quad (25b)$$

According to Theorem 3.1 the convergence rate of Eqs. (21) to Eqs. (25) scales with $d^{-(1/2+\kappa+\delta)} \log(d)$. Therefore the existence of the noisy ODEs above is circumscribed to the region

$$-1/2 < \kappa + \delta < 0 , \quad (26)$$

and presents a convergence trade-off absent in the other regimes: the faster one of the contributions of Eqs. (21) goes to zero, the worse is the convergence rate. In the present case, the more the learning term is attenuated, i.e. the more negative is $\kappa + \delta$, the worse the dynamics is described by Eqs. (25). Although the weights are updated, the correlation between the teacher and the student weights parametrized by the overlap matrix M remains fixed on its initial value M^0 , which is a fixed point of the dynamics under this scaling. Unsurprisingly, this leads to poor generalization capacity.

Red region (no ODEs) – If $\kappa + \delta < -1/2$, the stochastic process driven by the weight dynamics does not converge to deterministic ODEs under the assumptions of Theorem 3.1. We are then not able to state any claim about this regime.

Initialization and convergence – There are two additional features worth commenting on the high-dimensional dynamics and its connection to the mean-field/hydrodynamic approach, regarding initialization and the specialization transition.

In the ODE approach we discuss here, we always observe a first plateau where the teacher-student overlaps are all the same. This means all the hidden layer neurons learned the same linear separator. At this point, the two-layer network is essentially linear. This is called a *unspecialized* network in [16, 5]. In fact, this is a perfectly normal phenomenon, as with few samples even the Bayes-optimal solution would be unspecialized [24]. Only by running the dynamics long enough the student hidden neurons start to *specialize*, each of them learning a different sub-function so that the two-layer network can learn the non-trivial teacher.

Let us make two comments on this phenomenon: (i) while the "linear" learning in the unspecialized regime may remind the reader of the linear learning in the lazy regime [25, 26] of neural nets, the two phenomena are *completely* different. In lazy training, the learning is linear because weights change very little, so that the effective network is a linear approximation of the initial one. Here, instead, the weights are changing *considerably*, but each hidden neuron learns essentially the same function. (ii) If the ODEs are initialized with weights uncorrelated with the teacher, then the unspecialized regime is a fixed point of the ODEs: the student thus never specializes, at any time. Strikingly, such condition arises as well in the analysis of mean-field equations (see e.g. Theorem 2 in [27] that discusses the need to have *spread* initial conditions with a non-zero overlap with the teacher) to guarantee global convergence.

This raises the question about the precise dependence of the learning on the initialization condition in the high-dimensional regime, where a random start gets a vanishing ($1/\sqrt{d}$) overlap. This is a challenging problem that only recently has been studied (though in a simpler setting) in [28, 29, 30] who showed it yields an additional $\log(d)$ time-dependence. Generalizing these results for high-dimensional two-layer nets is an open question which we leave for future work.

4 Discussion, special cases, and simulations

To illustrate the phase diagram of Figure 1a, we present now several special cases for which we can perform simulations or numerically solve the set of ODEs. Henceforth, we take $\sigma(x) = \text{erf}(x/\sqrt{2})$, for which the expectations of the ODEs and of the population risk, Eq. (12), can be calculated analytically [5]. The explicit expressions are presented in Appendix C. Teacher weights are such that $\rho_{rs} = \delta_{rs}$. The initial student weights are chosen such that the dimension d can be varied without changing the initial conditions Q^0 , M^0 , P and consequently the initial population risk \mathcal{R}_0 . A detailed discussion can be found in Appendix D.

4.1 Saad & Solla scaling $\kappa = \delta = 0$

We start by recalling the well-known setting characterized by the point $\kappa = \delta = 0$. The convergence of the stochastic process for fixed learning rate and hidden layer width to Eqs. (23) was first obtained heuristically by [5]. In Figure 2 we recall this classical result by plotting the population risk dynamics for different noise levels. Dots represent simulations, while solid lines are obtained by integration of the ODEs, Eq. (23).

Learning is characterized by two phases after the initial decay. The first is the unspecialized plateau where all the teacher-student overlaps are approximately the same: $m_{jr} \approx m$. Waiting long enough, the dynamics reaches the *specialization* phase, where the student neurons start to *specialize*, i.e., their overlaps with one of the teacher neurons increase and consequently the population risk decreases. This specialization is discussed extensively in [5]. If $\Delta = 0$, the population risk goes asymptotically to zero. Instead, if $\Delta \neq 0$, the specialization phase presents a second plateau related to the noise Δ .

The asymptotic population risk \mathcal{R}_∞ related to the second plateau is proportional to $\gamma\Delta$ [6] in the high-dimensional limit $d \rightarrow \infty$ with p finite. As mentioned in the previous section, the expectation over $\mathcal{E}_j \mathcal{E}_l$ in Eq. (23a) prevents one from obtaining zero population risk for a noisy teacher.

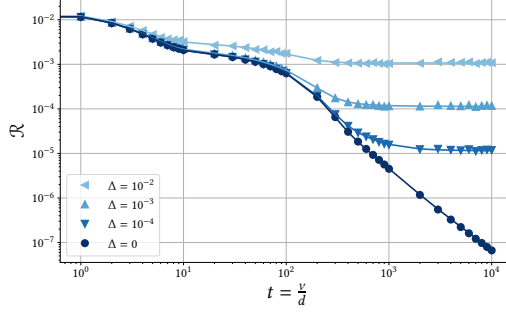
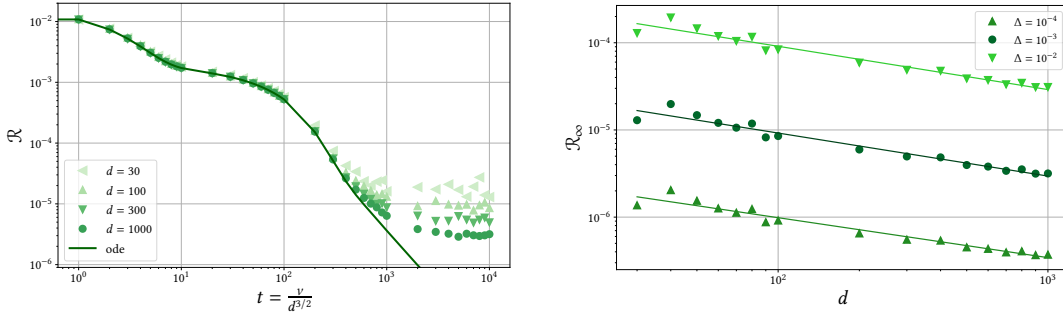


Figure 2: Population risk dynamics for $\kappa = \delta = 0$ (Saad & Solla scaling) : $p_0 = 8, k = 4, \rho_{rs} = \delta_{rs}$. Activation: $\sigma(x) = \text{erf}(x/\sqrt{2})$. Data distribution: $\mathbb{P}(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbf{1})$. Dots represent simulations ($d = 1000$), while solid lines are obtained by integration of the ODEs given by Eqs. (23).



(a) Population risk dynamics for $\kappa = 0$ and $\delta = 1/2$. Fixed noise $\Delta = 10^{-3}$ and varying d . Dots represent simulations, while the solid line is obtained by integration of the ODEs given by Eqs. (24). The data are compatible with the claim that as $d \rightarrow \infty$ the curve converges to zero population risk.

(b) Asymptotic population risk \mathcal{R}_∞ from simulations (dots) as a function of d for different noise levels under the scaling $\kappa = 0$ and $\delta = 1/2$. The fitted straight lines have slopes $-0.458, -0.494, -0.497$, for $\Delta = 10^{-4}, 10^{-3}, 10^{-2}$, respectively.

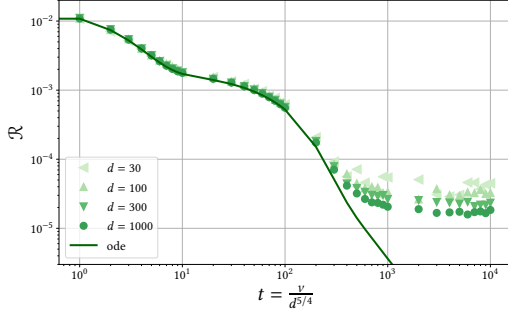
Figure 3: Network parameters: $p_0 = 8, k = 4, \rho_{rs} = \delta_{rs}$. Activation function: $\sigma(x) = \text{erf}(x/\sqrt{2})$. Data distribution: $\mathbb{P}(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbf{1})$.

4.2 Perfect learning for $\kappa = 0$

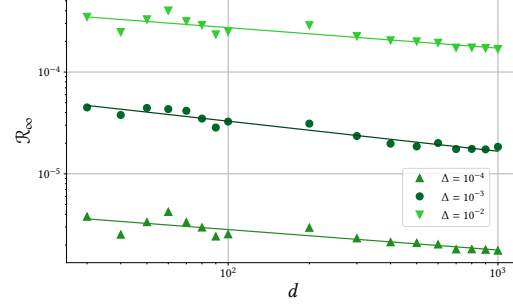
In this section we study the line $\kappa = 0$ with $\delta > 0$ of Figure 1a, for which Eqs. (24) with $\kappa = 0$ hold. We show that perfect learning can be asymptotically achieved in the realizable setting for any finite hidden layer width $p = p_0$. Keeping δ and Δ fixed, we have done simulations increasing the input layer dimension d . In Figure 3a we set $\delta = 1/2, \Delta = 10^{-3}$ and vary the input layer dimension. The bigger d is, the closer we are to the ODE-derived noiseless result.

Gathering the asymptotic population risk from simulations for varying d and Δ we perform a finite-size analysis to study the dependence of \mathcal{R}_∞ with d . This shows that the noise term goes to zero under this setting. In Figure 3b we plot \mathcal{R}_∞ versus d from simulations (dots) for different noise levels. We fit lines under the log-log scale showing that $\mathcal{R}_\infty \propto d^{-\delta}$, as expected. Figure 4 draws the same conclusion for $\delta = 1/4$.

As already stated, the interplay between the exponents directly affects the time scale. We end this subsection by graphically illustrating this fact through simulations. Setting the noise to $\Delta = 10^{-3}$ we

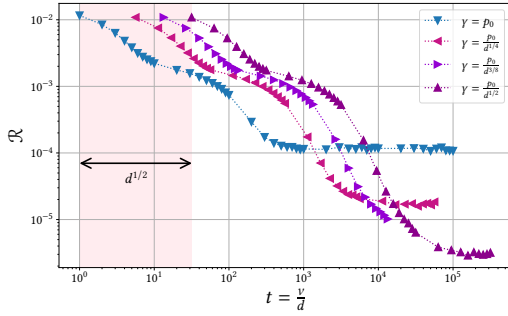


(a) Population risk dynamics for $\kappa = 0$ and $\delta = 1/4$. Fixed noise $\Delta = 10^{-3}$ and varying d . Dots represent simulations, while the solid line is obtained by integration of the ODEs given by Eqs. (24). The data are compatible with the claim that as $d \rightarrow \infty$ the curve converges to zero population risk.

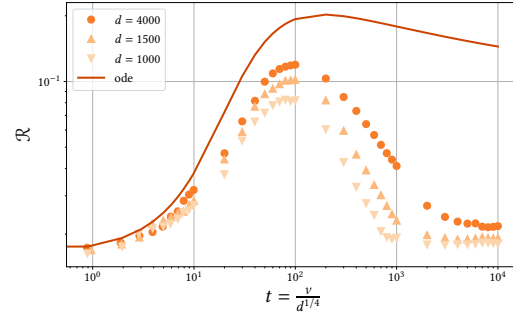


(b) Asymptotic population risk \mathcal{R}_∞ from simulations (dots) as a function of d for different noise levels under the scaling $\kappa = 0$ and $\delta = 1/4$. The fitted straight lines have slopes $-0.201, -0.295, -0.201$, for $\Delta = 10^{-4}, 10^{-3}, 10^{-2}$, respectively.

Figure 4: Network parameters: $p_0 = 8, k = 4, \rho_{rs} = \delta_{rs}$. Activation function: $\sigma(x) = \text{erf}(x/\sqrt{2})$. Data distribution: $\mathbb{P}(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbf{1})$.



(a) Simulations ($d = 1000$) for $\kappa = 0$ comparing different choices of the exponent δ . The final plateau is proportional to learning rate: $\mathcal{R}_\infty \propto \gamma \Delta$.



(b) Population risk dynamics for $\kappa = 0$ and $\delta = -3/8$. Dots represent simulations, while the solid line is obtained by integration of the ODEs given by Eqs. (25).

Figure 5: Network parameters $p_0 = 8, k = 4, \rho_{rs} = \delta_{rs}$. Noise level $\Delta = 10^{-3}$. Activation: $\sigma(x) = \text{erf}(x/\sqrt{2})$. Data distribution: $\mathbb{P}(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbf{1})$.

compare the cases $\delta = 0, 1/4, 3/8, 1/2$ in Figure 5a. All simulations are rendered on the scale $\delta t_0 = 1/d$ to illustrate the trade-off between asymptotic performance and training time.

4.3 Bad learning for $\kappa = 0$

We now quickly discuss the uncommon case of γ growing with d within the orange region. In Figure 5b we compare simulations varying d with the solution of the ODEs given by Eqs. (25). Both lead to poor results compared to the green and blue regions. Moreover, this regime presents strong finite-size effects, making it harder to observe the asymptotic ODEs at small sizes. However, the trend as d increases is very clear from the simulations. As discussed in Section 3, the more the learning term is attenuated on the ODEs, the worse they describe the dynamics.

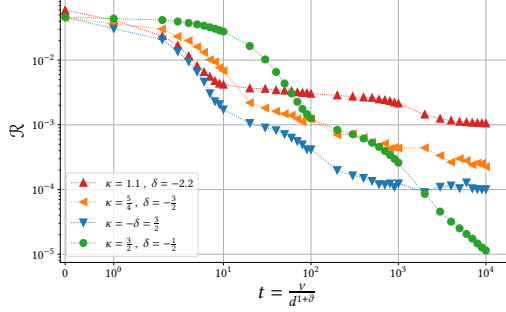


Figure 6: Simulations across different regions of Figure 1a. Networks parameters $d = 100$, $p = d^k$, $\gamma = d^{-\delta}$, $k = 4$, $\rho_{rs} = \delta_{rs}$. Noise: $\Delta = 10^{-3}$. Activation function: $\sigma(x) = \text{erf}(x/\sqrt{2})$. Data distribution: $\mathbb{P}(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbf{1})$. Time scaling: $\vartheta = \kappa + \delta$ for green and blue and $\vartheta = 2(\kappa + \delta)$ for orange. The colors match Figure 1a.

4.4 Large hidden layer: $\kappa > 0$

Finishing our voyage through Figure 1a with examples, we briefly discuss the case where both input and hidden layer widths are large. Although Theorem 3.1 provides non-asymptotic guarantees for $\kappa > 0$, the number of coupled ODEs grows quadratically with p , making the task of solving them rather challenging. Thus, we present simulations that illustrate the regions of Figure 1a. Fixing $d = 100$ we show in Figure 6 learning curves for different values of κ and δ . The colors are chosen to match their respective regions in the phase diagram.

Due to the relatively small sizes used in Figure 6, the green dots seem to decrease towards perfect learning, even when $\delta < 0$, provided that κ is large enough, as is predicted by the phase diagram in Figure 1a. Moreover, since d is not large enough, when the parameters are within the orange region the finite-size effects actually dominates, similarly to Figure 5b. The learning contribution still plays a role and the asymptotic population risk is similar to the case $\kappa = \delta = 0$. Within the red region, which is out of scope of our theory, the simulation gets stuck on a plateau with larger population risk.

5 Conclusion

Building up on classical statistical physics approaches and extending them to a broad range of learning rate, time scales, and hidden layer width, we rendered a sharp characterisation of the performance of SGD for two-layer neural networks in high-dimensions. Our phase diagram describes the possible learning scenarios, characterizing learning regimes which had not been addressed by previous classical works using ODEs. Crucially, our key conclusions do not rely on an explicit solution, as our theory allows the characterization of the learning dynamics *without* solving the system of ODEs. The introduction of scaling factors is non-trivial and has deep implications. Our generalized description enlightens the trade-off between learning rate and hidden layer width, which has also been crucial in the mean-field theories.

Acknowledgements

We thank Gérard Ben Arous, Lenaïc Chizat, Maria Refinetti and Sebastian Goldt for discussions. We acknowledge funding from the ERC under the European Union’s Horizon 2020 Research and Innovation Pro-

gram Grant Agreement 714608- SMiLe. RV was partially financed by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Brasil (CAPES) - Finance Code 001. RV is grateful to EPFL and IdePHICS lab for their generous hospitality during the realization of this project.

Appendix

A Deterministic scaling limit of stochastic processes

In order to show the deterministic scaling of online SGD under a proper chosen time scale, we will make use of a convergence result by [21, 31], which is adapted below in Theorem A.1.

Theorem A.1 (Deterministic scaling limit of stochastic processes). *Consider a d -dimension discrete time stochastic process sequence, $\{\Omega^v ; v = 0, 1, 2, \dots, [S\tau]\}_{S=1,2,\dots}$ for some $\tau > 0$. The increment $\Omega^{v+1} - \Omega^v$ is assumed to be decomposable into three parts,*

$$\Omega^{v+1} - \Omega^v = \frac{1}{S}\psi(\Omega^v) + \Lambda^v + \Gamma^v, \quad (\text{A.1})$$

such that

Assumption A.1.1. The process $\tilde{\Lambda}^v \equiv \sum_{v'=0}^v \Lambda^{v'}$ is a martingale and $\mathbb{E}\|\Lambda^v\|^2 \leq C(\tau)^2/S^{1+\epsilon_1}$ for some $\epsilon_1 > 0$.

Assumption A.1.2. $\mathbb{E}\|\Gamma^v\| \leq C(\tau)/S^{1+\epsilon_2}$ for some $\epsilon_2 > 0$.

Assumption A.1.3. The function $\psi(\Omega)$ is Lipschitz, i.e, $\|\psi(\Omega) - \psi(\tilde{\Omega})\| \leq C\|\Omega - \tilde{\Omega}\|$ for any Ω and $\tilde{\Omega}$.

Let $\Omega(t)$, with $0 \leq t \leq \tau$, be a continuous stochastic process such that $\Omega(t) = \Omega^v$ with $v = [St]$. Define the deterministic ODE

$$\frac{d}{dt}\bar{\Omega}(t) = \psi(\bar{\Omega}(t)), \quad (\text{A.2})$$

with $\bar{\Omega}(0) = \bar{\Omega}_0$.

Then, if assumptions A.1.1 to A.1.3 hold and assuming $\mathbb{E}\|\Omega^0 - \bar{\Omega}_0\| < C/S^{\epsilon_3}$ for some $\epsilon_3 > 0$ then we have for any finite S :

$$\mathbb{E}\left\|\Omega^v - \bar{\Omega}\left(\frac{v}{S}\right)\right\| \leq C(\tau)e^{c\tau}S^{-\min\{\frac{1}{2}\epsilon_1, \epsilon_2, \epsilon_3\}}, \quad (\text{A.3})$$

where $\bar{\Omega}(\cdot)$ is the solution of Eq.(A.2).

Proof. The reader interested in the proof is referred to the supplementary materials of [21, 31]. \square

Although the theorem wasn't originally proven in the $p \rightarrow \infty$ setting, a glance at its proof shows that it still holds upon replacing $C(\tau)$ by $C(p, \tau)$ in Assumption A.1.1 and A.1.2, as well as Equation (A.3). We choose $\|\cdot\|$ to be the L^∞ norm, since it suits better the $p \rightarrow \infty$ scaling. The S in Theorem A.1 corresponds to $1/\delta t$, where δt is defined in Theorem 3.1.

Following [21], we define for $j, l \in [p]$

$$\Psi_{jl}(\Omega; \mathbf{x}) = \frac{Y}{pd\delta t} \left(\mathcal{E}_j^v \lambda_l^v + \mathcal{E}_l^v \lambda_j^v \right) + \frac{Y^2}{p^2 d \delta t} \mathcal{E}_j^v \mathcal{E}_l^v,$$

and

$$\psi_{jl}(\Omega) = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbb{I})} \left[\Psi_{jl}(\Omega; \mathbf{x}) \right].$$

The functions Ψ, ψ are similarly defined on $[p] \times [p+1, p+k]$. With that, we write

$$\Omega^{\nu+1} - \Omega^\nu = \frac{1}{S}\psi(\Omega) + \underbrace{\frac{1}{S}(\Psi(\Omega^\nu; \mathbf{x}) - \psi(\Omega^\nu))}_{\Lambda^\nu} + \Gamma^\nu,$$

where for $j, l \in [p]$

$$\Gamma_{jl}^\nu = \frac{Y^2}{p^2 d^2} (\|\mathbf{x}\|_2^2 - d) \mathcal{E}_j^\nu \mathcal{E}_l^\nu.$$

The main obstacle to bounding Λ^ν and Γ^ν is the fact that the q_{jj} can a priori diverge to infinity. Our first task is therefore to show that this does not happen; as a proxy we show a subgaussian-like moment bound:

$$\mathbb{E}[(q_{jj}^\nu)^t] \leq \left(C(\tau) + \frac{ct}{S}\right)^t.$$

Equipped with the above bound, controlling $\mathbb{E}\|\Lambda^\nu\|^2$ and $\mathbb{E}\|\Gamma^\nu\|$ becomes fairly easy. All proof details are in the below sections.

A.1 Preliminaries: bounding the q_{jj}

Since σ is L -Lipschitz, we have by the Cauchy-Schwarz inequality

$$(\mathcal{E}^\nu)^2 \leq \frac{3L^2}{k} \sum_{r=1}^k (\lambda_r^*)^2 + \frac{3L^2}{p} \sum_{j=1}^p (\lambda_j)^2 + 3\Delta\zeta^2 \equiv \Phi^\nu \quad (\text{A.4})$$

Define

$$s^\nu = \mathbb{E}\Phi^\nu = \frac{3L^2}{k} \sum_{r=1}^k \rho_{rr} + \frac{3L^2}{p} \sum_{j=1}^p q_{jj}^\nu + 3\Delta$$

Assumption 1 in Theorem 3.1 implies that

$$|q_{jj}^{\nu+1} - q_{jj}^\nu| \leq \frac{1}{S} \left(c_1 (\lambda_j^\nu)^2 + c_2 (\mathcal{E}^\nu)^2 \right)$$

where c_1, c_2 are absolute constants. Summing those inequalities yield

$$|s_{\nu+1} - s^\nu| \leq \frac{c_3}{S} \Phi^\nu,$$

and finally

$$\mathbb{E}_\nu[s^{\nu+1}] \leq s^\nu \left(1 + \frac{c_3}{S} \right) \leq s^\nu e^{c_3/S}.$$

As a result, we have for any $0 \leq \nu \leq S\tau$

$$\mathbb{E}[s^\nu] \leq c_4 e^{c_3\tau}. \quad (\text{A.5})$$

For simplicity, let q^ν denote any of the q_{jj}^ν . We have, for all $t \geq 0$,

$$(q^{\nu+1})^t - (q^\nu)^t = t(q^\nu)^{t-1}(q^{\nu+1} - q^\nu) + O\left(\frac{t^2}{S^2}\right),$$

where the remainder term has bounded expectation. Again, we write

$$|(q^{v+1})^t - (q^v)^t| \leq t(q^v)^{t-1} \frac{1}{S} (c_1(\mathcal{E}^v)^2 + c_2(\lambda_i^v)^2) + \frac{c_5 t^2}{S^2}.$$

By Assumption 3, the q_{ii}^v are bounded from below by a constant, hence

$$\mathbb{E}_v[(q^{v+1})^t] \leq (q^v)^t \left(1 + \frac{c_6 t}{S}\right) + O\left(\frac{c_5 t^2}{S^2}\right)$$

This implies that for any $t \geq 0$ and $0 \leq v \leq S\tau$,

$$\mathbb{E}[(q^v)^t] \leq \left(c_7 + \frac{c_5 t^2}{S}\right) e^{c_6 \tau} \leq \left(C(\tau) + \frac{c_5 t}{S}\right)^t \quad (\text{A.6})$$

A.2 Assumption A.1.1

We have for all $i, j \in [p+k]$,

$$\left(\Omega_{ij}^{v+1} - \mathbb{E}_v[\Omega_{ij}^{v+1}]\right)^2 \leq 2 \left((\Omega_{ij}^{v+1} - \Omega_{ij}^v)^2 + (\Omega_{ij}^v - \mathbb{E}_v[\Omega_{ij}^{v+1}])^2 \right).$$

As a consequence,

$$\mathbb{E}\|\Lambda^v\|^2 \leq 4 \max_{i,j} (\Omega_{ij}^{v+1} - \Omega_{ij}^v)^2.$$

Now, by definition,

$$(q_{ij}^{v+1} - q_{ij}^v)^2 \leq \frac{L}{S^2} (c_1(\mathcal{E}^v)^2 + c_2|\mathcal{E}^v|(|\lambda_i| + |\lambda_j|))^2 \leq \frac{L}{S^2} (c_3(\mathcal{E}^v)^4 + c_4(\max_{\ell} \lambda_{\ell}^v)^4),$$

The term in $(\mathcal{E}^v)^4$ is bounded by the same techniques as the last section. For the second term,

$$\mathbb{E}_v \left[(\max_{\ell} \lambda_{\ell}^v)^4 \right] \leq c_5 \log(p)^2 \left(\max_{\ell} q_{\ell\ell}^v \right)^4,$$

and we can write for any $t \geq 0$

$$\max_{\ell} (q_{\ell\ell}^v)^4 \leq \left(\sum_{\ell} (q_{\ell\ell}^v)^t \right)^{4/t}.$$

By Jensen's inequality, for $t \geq 4$

$$\mathbb{E} \left[\left(\max_{\ell} q_{\ell\ell}^v \right)^4 \right] \leq \left(\sum_{\ell} \mathbb{E}[(q_{\ell\ell}^v)^t] \right)^{4/t} \leq p^{4/t} \left(C(\tau) + \frac{c_6 t}{S} \right)^4,$$

using (A.6). Choosing $t = 4 \log(p) \ll S$ shows that

$$\mathbb{E} \left[\max_{i,j} (q_{ij}^{v+1} - q_{ij}^v)^2 \right] \leq \frac{C(\tau) \log(p)^2}{S^2}$$

A similar bound holds for the m_{ij} , and hence

$$\mathbb{E}\|\Lambda^v\|^2 \leq \frac{c_5 \log(p)^2}{S^2},$$

which implies Assumption A.1.1 with $\epsilon_1 = 1$ and $C(p, \tau) = C'(\tau) \log(p)$.

A.3 Assumption A.1.2

Since σ is Lipschitz, for any $i, j \in [p]$

$$\mathcal{E}_i^v \mathcal{E}_j^v \leq L^2 (\mathcal{E}^v)^2.$$

Hence,

$$\begin{aligned} \mathbb{E}[\|\Gamma^v\|_\infty] &\leq \frac{L^2 \gamma^2}{d^2 p^2} \mathbb{E}[(\|\mathbf{x}\|_2^2 - d) \Phi^v] \\ &\leq \frac{L^2 \gamma^2}{d^2 p^2} \left(\frac{1}{2\sqrt{d}} \mathbb{E}[(\|\mathbf{x}\|_2^2 - d)^2] + \frac{\sqrt{d}}{2} \mathbb{E}[(\mathcal{E}^v)^4] \right). \end{aligned}$$

The first expectation is the variance of a χ_d^2 random variable, which is equal to $2d$, and the second expectation is bounded by the same methods as the above sections. The term in brackets is therefore bounded by $c_1 \sqrt{d}$, and

$$\mathbb{E}[\|\Gamma^v\|_\infty] \leq c_2 \frac{\gamma^2}{d^{3/2} p^2}$$

Finally, since for any $y > 0$ we have $y^2 \leq \max(y, y^2)^{3/2}$, letting $y = \gamma/p$ we find

$$\mathbb{E}[\|\Gamma^v\|_\infty] \leq c_2 \max\left(\frac{\gamma}{pd}, \frac{\gamma^2}{p^2 d}\right)^{3/2} \leq c_3 (\delta t)^{3/2},$$

hence Assumption A.1.2 is true with $\epsilon_2 = 1/2$.

A.4 $\sqrt{\cdot}$ -Lipschitz property

Let $\Omega, \Omega' \in \mathbb{R}^{(p+k) \times (p+k)}$, we can write the (i, j) coefficient of $\psi(\Omega)$ as $f_{ij}(\sqrt{\Omega})$, where

$$\begin{aligned} f &: \mathbb{R}^{(p+k) \times (p+k)} \rightarrow \mathbb{R} \\ A &\mapsto \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, I_{p+k})} [g_{ij}(Ax)] \end{aligned}$$

The same arguments as above show that the function f is Lipschitz, and hence for some constant L'' we have

$$\|\psi(\Omega) - \psi(\Omega')\| \leq L'' \|\sqrt{\Omega} - \sqrt{\Omega'}\|.$$

B A lemma on ODE perturbation

In this section, we prove a proposition that bounds the difference between an ODE solution and a perturbed version, for a bounded time t .

Theorem B.1. *Let $f, g : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be two L -Lipschitz functions, and consider the following differential equations in \mathbb{R}^n :*

$$\begin{aligned} \frac{dx}{dt} &= f(x) + \epsilon g(x), \\ \frac{dy}{dt} &= f(y), \end{aligned}$$

where $\epsilon > 0$, and with the initial condition $\mathbf{x}(0) = \mathbf{y}(0)$. Then, if $\tau > 0$ is fixed, we have

$$\|\mathbf{x}(t) - \mathbf{y}(t)\|_2 \leq c\epsilon e^{L\tau}$$

for any $0 \leq t \leq \tau$, with c a constant independent from ϵ, τ .

Before proving this proposition, we begin with a small lemma:

Lemma B.2. Let $a, b > 0$, and $z : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ a function satisfying

$$\frac{dz}{dt} = az + b\sqrt{z}$$

with $z(0) = 0$. Then, for some constant $c > 0$, we have

$$z(t) \leq c \frac{b^2 e^{at}}{a^2} \quad \text{for all } t \geq 0$$

Proof. Upon considering the function $a^2 z(t/a)/b^2$ instead, we can assume that $a = b = 1$. Then, we have

$$\frac{dz}{dt} \leq \max(z, 1) + \max(\sqrt{z}, 1),$$

and the RHS is an increasing function. Hence, if \tilde{z} is a solution of

$$\frac{d\tilde{z}}{dt} = \max(z, 1) + \max(\sqrt{\tilde{z}}, 1),$$

with $\tilde{z}(0) = 0$, then $z(t) \leq \tilde{z}(t)$ for all $t \geq 0$. Since the RHS of the above equation is Lipschitz everywhere, we can apply the Picard–Lindelöf theorem, and check that the unique solution to this equation is

$$\tilde{z}(t) = \begin{cases} 2t & \text{if } t \leq \frac{1}{2} \\ (c_1 e^t - c_2)^2 & \text{otherwise} \end{cases},$$

where c_1 and c_2 are ad hoc constants. The lemma then follows from adjusting the constant c as needed. \square

We are now in a position to show Theorem B.1:

Proof. Assume for simplicity that $\mathbf{x}(0) = \mathbf{y}(0) = \mathbf{0}$. We begin by bounding $\mathbf{x}(t)$; we have

$$\frac{d\|\mathbf{x}\|^2}{dt} = 2\mathbf{x}^\top \frac{d\mathbf{x}}{dt} \leq 2\|\mathbf{x}\| \|f(\mathbf{x}) + \epsilon g(\mathbf{x})\|.$$

By the Lipschitz condition,

$$\|f(\mathbf{x}) + \epsilon g(\mathbf{x})\| \leq \|f(\mathbf{0}) + \epsilon g(\mathbf{0})\| + \frac{L}{2}\|\mathbf{x}\|,$$

so that

$$\frac{d\|\mathbf{x}\|^2}{dt} \leq L\|\mathbf{x}\|^2 + 2\|f(\mathbf{0}) + \epsilon g(\mathbf{0})\| \|\mathbf{x}\|.$$

Applying Lemma B.2 and taking square roots on each side,

$$\|\mathbf{x}(t)\| \leq c \frac{\|f(\mathbf{0}) + \epsilon g(\mathbf{0})\|}{L} e^{Lt/2} \leq c \frac{\|f(\mathbf{0}) + \epsilon g(\mathbf{0})\|}{L} e^{L\tau/2}, \quad (\text{B.1})$$

for any $0 \leq t \leq \tau$. Now, similarly,

$$\begin{aligned} \frac{d\|\mathbf{x} - \mathbf{y}\|^2}{dt} &\leq 2\|\mathbf{x} - \mathbf{y}\| \left\| \frac{d(\mathbf{x} - \mathbf{y})}{dt} \right\| \\ &\leq 2\|\mathbf{x} - \mathbf{y}\| \|f(\mathbf{x}) - f(\mathbf{y}) + \epsilon g(\mathbf{x})\| \\ &\leq L\|\mathbf{x} - \mathbf{y}\|^2 + 2\epsilon\|g(\mathbf{x})\| \|\mathbf{x} - \mathbf{y}\| \\ &\leq L\|\mathbf{x} - \mathbf{y}\|^2 + \epsilon \left(\|g(\mathbf{0})\| + c\|f(\mathbf{0}) + \epsilon g(\mathbf{0})\| e^{L\tau/2} \right) \|\mathbf{x} - \mathbf{y}\|, \end{aligned}$$

having used (B.1) on the last line. This is again the setting of Lemma B.2, which gives

$$\|\mathbf{x} - \mathbf{y}\| \leq c_1 \epsilon e^{L\tau/2} \frac{e^{L\tau/2}}{L} \leq c_2 \epsilon e^{L\tau}.$$

□

C Expectations over the local fields

In this appendix we present the explicit expressions from the expectations of the local fields used to compute the population risk and the ODE terms.

C.1 Population risk

We write the population risk (12) as

$$\begin{aligned} \mathcal{R}(\Omega) &= \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | \mathbf{0}, \Omega)} \mathbb{E}_{\zeta \sim \mathbb{P}(\zeta)} \left[\left(\hat{f}(\lambda) - f(\lambda^*) \right)^2 \right] \\ &= \mathcal{R}_t(\mathbf{P}) + \mathcal{R}_s(\mathbf{Q}) + \mathcal{R}_{\text{st}}(\mathbf{P}, \mathbf{Q}, \mathbf{M}), \end{aligned} \quad (\text{C.1})$$

with

$$\mathcal{R}_t \equiv \mathbb{E}_{\lambda^* \sim \mathcal{N}(\lambda^* | \mathbf{0}, \mathbf{P})} \left[f(\lambda^*)^2 \right] = \frac{1}{k^2} \sum_{r,s=1}^k \mathbb{E}_{\lambda^* \sim \mathcal{N}(\lambda^* | \mathbf{0}, \mathbf{P})} \left[\sigma(\lambda_r^*) \sigma(\lambda_s^*) \right] \quad (\text{C.2a})$$

$$\mathcal{R}_s \equiv \mathbb{E}_{\lambda \sim \mathcal{N}(\lambda | \mathbf{0}, \mathbf{Q})} \left[\hat{f}(\lambda)^2 \right] = \frac{1}{p^2} \sum_{j,l=1}^k \mathbb{E}_{\lambda \sim \mathcal{N}(\lambda | \mathbf{0}, \mathbf{Q})} \left[\sigma(\lambda_j) \sigma(\lambda_l) \right], \quad (\text{C.2b})$$

$$\mathcal{R}_{\text{st}} \equiv \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | \mathbf{0}, \Omega)} \left[\hat{f}(\lambda) f(\lambda^*) \right] = -\frac{2}{pk} \sum_{j=1}^p \sum_{r=1}^k \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | \mathbf{0}, \Omega)} \left[\sigma(\lambda_j) \sigma(\lambda_r^*) \right] \quad (\text{C.2c})$$

Define the vector $\boldsymbol{\lambda}^{\alpha\beta} \equiv (\lambda^\alpha, \lambda^\beta)^\top \in \mathbb{R}^2$, where the upper indices on the components indicate they may refer to student or teacher local fields. Consider the covariance matrix on the subspace spanned by $\boldsymbol{\lambda}^{\alpha\beta}$:

$$\Omega^{\alpha\beta} \equiv \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | \mathbf{0}, \Omega)} \left[\boldsymbol{\lambda}^{\alpha\beta} \left(\boldsymbol{\lambda}^{\alpha\beta} \right)^\top \right] \in \mathbb{R}^{2 \times 2}. \quad (\text{C.3})$$

For $\sigma(x) = \text{erf}(x/\sqrt{2})$ the expectations in Eqs. (C.2) are in general given by [5]

$$\mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \left[\sigma(\lambda^\alpha) \sigma(\lambda^\beta) \right] = \frac{1}{\pi} \arcsin \left(\frac{\Omega_{12}^{\alpha\beta}}{\sqrt{(1 + \Omega_{11}^{\alpha\beta})(1 + \Omega_{22}^{\alpha\beta})}} \right). \quad (\text{C.4})$$

where $\Omega_{jl}^{\alpha\beta} \equiv (\Omega^{\alpha\beta})_{jl}$ is an element of the covariance matrix given by Eq. (C.3).

Explicitly, the population risk contributions are

$$\mathcal{R}_t(\mathbf{P}) = \frac{1}{k^2} \sum_{r,s=1}^k \frac{1}{\pi} \arcsin \left(\frac{\rho_{rs}}{\sqrt{(1 + \rho_{rr})(1 + \rho_{ss})}} \right), \quad (\text{C.5a})$$

$$\mathcal{R}_s(\mathbf{Q}) = \frac{1}{p^2} \sum_{j,l=1}^k \frac{1}{\pi} \arcsin \left(\frac{q_{jl}}{\sqrt{(1 + q_{jj})(1 + q_{ll})}} \right), \quad (\text{C.5b})$$

$$\mathcal{R}_{st}(\mathbf{P}, \mathbf{Q}, \mathbf{M}) = -\frac{2}{pk} \sum_{j=1}^p \sum_{r=1}^k \frac{1}{\pi} \arcsin \left(\frac{m_{jr}}{\sqrt{(1 + q_{jj})(1 + \rho_{rr})}} \right). \quad (\text{C.5c})$$

C.2 ODE contributions

From the update equations, we first consider the expectations linear in \mathcal{E}_j :

$$\begin{aligned} \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \mathbb{E}_{\zeta \sim \mathbb{P}(\zeta)} \left[\mathcal{E}_j \lambda_l \right] &= \frac{1}{k} \sum_{r'=1}^k \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \left[\sigma'(\lambda_j) \lambda_l \sigma(\lambda_{r'}) \right] \\ &\quad - \frac{1}{p} \sum_{l'=1}^p \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \left[\sigma'(\lambda_j) \lambda_l \sigma(\lambda_{l'}) \right], \end{aligned} \quad (\text{C.6a})$$

$$\begin{aligned} \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \mathbb{E}_{\zeta \sim \mathbb{P}(\zeta)} \left[\mathcal{E}_j \lambda_r^* \right] &= \frac{1}{k} \sum_{r'=1}^k \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \left[\sigma'(\lambda_j) \lambda_r^* \sigma(\lambda_{r'}) \right] \\ &\quad - \frac{1}{p} \sum_{l'=1}^p \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \left[\sigma'(\lambda_j) \lambda_r^* \sigma(\lambda_{l'}) \right]. \end{aligned} \quad (\text{C.6b})$$

Define the vector $\boldsymbol{\lambda}^{\alpha\beta\gamma} \equiv (\lambda^\alpha, \lambda^\beta, \lambda^\gamma)^\top \in \mathbb{R}^3$, where the upper indices on the components indicate they may refer to student or teacher local fields. Consider the covariance matrix on the subspace spanned by $\boldsymbol{\lambda}^{\alpha\beta\gamma}$:

$$\Omega^{\alpha\beta\gamma} \equiv \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \left[\boldsymbol{\lambda}^{\alpha\beta\gamma} \left(\boldsymbol{\lambda}^{\alpha\beta\gamma} \right)^\top \right] \in \mathbb{R}^{3 \times 3}. \quad (\text{C.7})$$

For $\sigma(x) = \text{erf}(x/\sqrt{2})$ the expectations in Eqs. (C.6) are given by [5]

$$\mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \left[\sigma'(\lambda^\alpha) \lambda^\beta \sigma(\lambda^\gamma) \right] = \frac{2}{\pi} \frac{\Omega_{23}^{\alpha\beta\gamma} (1 + \Omega_{11}^{\alpha\beta\gamma}) - \Omega_{12}^{\alpha\beta\gamma} \Omega_{13}^{\alpha\beta\gamma}}{(1 + \Omega_{11}^{\alpha\beta\gamma}) \sqrt{(1 + \Omega_{11}^{\alpha\beta\gamma})(1 + \Omega_{33}^{\alpha\beta\gamma}) - (\Omega_{13}^{\alpha\beta\gamma})^2}}, \quad (\text{C.8})$$

where $\Omega_{jl}^{\alpha\beta\gamma} \equiv (\Omega^{\alpha\beta\gamma})_{jl}$ is an element of the covariance matrix given by Eq. (C.7). As examples, we write explicitly:

$$\Omega^{jlr'} = \begin{bmatrix} q_{jj} & q_{jl} & m_{jr'} \\ q_{jl} & q_{ll} & m_{lr'} \\ m_{jr'} & m_{lr'} & \rho_{r'r'} \end{bmatrix}, \quad \Omega^{jrr'} = \begin{bmatrix} q_{jj} & m_{jr} & m_{jr'} \\ m_{jr} & \rho_{rr} & \rho_{rr'} \\ m_{jr'} & \rho_{rr} & \rho_{r'r'} \end{bmatrix}. \quad (\text{C.9})$$

The quadratic contribution in \mathcal{E}_j is given by

$$\begin{aligned} \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} \mathbb{E}_{\zeta \sim \mathcal{P}(\zeta)} [\mathcal{E}_j \mathcal{E}_l] &= \frac{1}{k^2} \sum_{r, r'=1}^k \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} [\sigma'(\lambda_j) \sigma'(\lambda_l) \sigma(\lambda_r^*) \sigma(\lambda_{r'}^*)] \\ &+ \frac{1}{p^2} \sum_{j', l'=1}^p \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} [\sigma'(\lambda_j) \sigma'(\lambda_l) \sigma(\lambda_{j'}) \sigma(\lambda_{l'})] \\ &- \frac{2}{pk} \sum_{l'=1}^p \sum_{r=1}^k \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} [\sigma'(\lambda_j) \sigma'(\lambda_l) \sigma(\lambda_r^*) \sigma(\lambda_{l'})] \\ &+ \Delta \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} [\sigma'(\lambda_j) \sigma'(\lambda_l)] \end{aligned} \quad (\text{C.10})$$

The solution of the noise-dependent term can be constructed with the covariance matrix (C.3) and is given by [6]

$$\mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} [\sigma'(\lambda^\alpha) \sigma'(\lambda^\beta)] = \frac{2}{\pi} \frac{1}{\sqrt{1 + \Omega_{11}^{\alpha\beta} + \Omega_{22}^{\alpha\beta} + \Omega_{11}^{\alpha\beta} \Omega_{22}^{\alpha\beta} - (\Omega_{12}^{\alpha\beta})^2}} \quad (\text{C.11})$$

Similarly, one can define the vector $\lambda^{\alpha\beta\gamma\delta} \equiv (\lambda^\alpha, \lambda^\beta, \lambda^\gamma, \lambda^\delta)^\top \in \mathbb{R}^4$ and write the covariance matrix on the subspace spanned by $\lambda^{\alpha\beta\gamma\delta}$:

$$\Omega^{\alpha\beta\gamma\delta} \equiv \mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} [\lambda^{\alpha\beta\gamma\delta} (\lambda^{\alpha\beta\gamma\delta})^\top] \in \mathbb{R}^{4 \times 4}. \quad (\text{C.12})$$

For $\sigma(x) = \text{erf}(x/\sqrt{2})$ the expectations in Eqs. (C.10) are given by [5]

$$\mathbb{E}_{\lambda, \lambda^* \sim \mathcal{N}(\lambda, \lambda^* | 0, \Omega)} [\sigma'(\lambda^\alpha) \sigma'(\lambda^\beta) \sigma(\lambda^\gamma) \sigma(\lambda^\delta)] = \frac{4}{\pi^2} \frac{1}{\sqrt{\bar{\Omega}_0^{\alpha\beta\gamma\delta}}} \arcsin \left(\frac{\bar{\Omega}_1^{\alpha\beta\gamma\delta}}{\sqrt{\bar{\Omega}_2^{\alpha\beta\gamma\delta} \bar{\Omega}_3^{\alpha\beta\gamma\delta}}} \right), \quad (\text{C.13})$$

with

$$\bar{\Omega}_0^{\alpha\beta\gamma\delta} \equiv \left(1 + \Omega_{11}^{\alpha\beta\gamma\delta}\right) \left(1 + \Omega_{22}^{\alpha\beta\gamma\delta}\right) - \left(\Omega_{12}^{\alpha\beta\gamma\delta}\right)^2, \quad (\text{C.14a})$$

$$\begin{aligned} \bar{\Omega}_1^{\alpha\beta\gamma\delta} &\equiv \bar{\Omega}_0^{\alpha\beta\gamma\delta} \Omega_{34}^{\alpha\beta\gamma\delta} - \Omega_{23}^{\alpha\beta\gamma\delta} \Omega_{24}^{\alpha\beta\gamma\delta} \left(1 + \Omega_{11}^{\alpha\beta\gamma\delta}\right) - \Omega_{13}^{\alpha\beta\gamma\delta} \Omega_{14}^{\alpha\beta\gamma\delta} \left(1 + \Omega_{22}^{\alpha\beta\gamma\delta}\right) \\ &+ \Omega_{12}^{\alpha\beta\gamma\delta} \Omega_{13}^{\alpha\beta\gamma\delta} \Omega_{24}^{\alpha\beta\gamma\delta} + \Omega_{12}^{\alpha\beta\gamma\delta} \Omega_{14}^{\alpha\beta\gamma\delta} \Omega_{23}^{\alpha\beta\gamma\delta}, \end{aligned} \quad (\text{C.14b})$$

$$\begin{aligned} \bar{\Omega}_2^{\alpha\beta\gamma\delta} &\equiv \bar{\Omega}_0^{\alpha\beta\gamma\delta} \left(1 + \Omega_{44}^{\alpha\beta\gamma\delta}\right) - \left(\Omega_{24}^{\alpha\beta\gamma\delta}\right)^2 \left(1 + \Omega_{11}^{\alpha\beta\gamma\delta}\right) - \left(\Omega_{13}^{\alpha\beta\gamma\delta}\right)^2 \left(1 + \Omega_{22}^{\alpha\beta\gamma\delta}\right) \\ &+ 2\Omega_{12}^{\alpha\beta\gamma\delta} \Omega_{13}^{\alpha\beta\gamma\delta} \Omega_{23}^{\alpha\beta\gamma\delta}, \end{aligned} \quad (\text{C.14c})$$

$$\begin{aligned} \bar{\Omega}_3^{\alpha\beta\gamma\delta} &\equiv \bar{\Omega}_0^{\alpha\beta\gamma\delta} \left(1 + \Omega_{44}^{\alpha\beta\gamma\delta}\right) - \left(\Omega_{24}^{\alpha\beta\gamma\delta}\right)^2 \left(1 + \Omega_{11}^{\alpha\beta\gamma\delta}\right) - \left(\Omega_{14}^{\alpha\beta\gamma\delta}\right)^2 \left(1 + \Omega_{22}^{\alpha\beta\gamma\delta}\right) \\ &+ 2\Omega_{12}^{\alpha\beta\gamma\delta} \Omega_{14}^{\alpha\beta\gamma\delta} \Omega_{24}^{\alpha\beta\gamma\delta}. \end{aligned} \quad (\text{C.14d})$$

C.3 From gradient flow to local fields

Consider the gradient flow approximation

$$\begin{aligned}\frac{d\mathbf{w}_j}{dt} &= -\nabla_{\mathbf{w}_j} \mathcal{R}(\mathbf{W}, \mathbf{W}^*) \\ &= -\frac{1}{p\sqrt{d}} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbb{1})} [\mathbf{x} \sigma'(\lambda_j) \mathcal{E}].\end{aligned}$$

Now, since for any $\mathbf{x}^\top \mathbf{y}$, we have

$$\frac{d(\mathbf{x}^\top \mathbf{y})}{dt} = \mathbf{x}^\top \frac{d\mathbf{y}}{dt} + \mathbf{y}^\top \frac{d\mathbf{x}}{dt},$$

we find

$$\frac{dq_{jl}}{dt} = -\frac{1}{pd} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbb{1})} [(\sigma'(\lambda_j)\lambda_l + \sigma'(\lambda_l)\lambda_j) \mathcal{E}].$$

Recalling the definition $\mathcal{E}_j = \sigma'(\lambda_j) \mathcal{E}$, the terms present inside the expectation are exactly those in the learning term of Eq.(14).

D Initial conditions and symmetric teacher

In this work we have constructed teacher matrices $\mathbf{W}^* \in \mathbb{R}^{k \times d}$ in order to have

$$\rho_{rs} = \frac{\mathbf{w}_r^{*\top} \mathbf{w}_s^*}{d} = \delta_{rs}, \quad (\text{D.1})$$

where $\mathbf{w}_r^* \equiv [\mathbf{W}^*]_r \in \mathbb{R}^d$ is the r -th row of the matrix \mathbf{W}^* . We have started by sampling k vectors of dimension d uniformly on a ball of radius \sqrt{d} . Then we constructed an orthonormal basis using singular value decomposition.

The initial student weights $\mathbf{W}^0 \in \mathbb{R}^{p \times d}$ were taken as

$$\mathbf{W}^0 = \mathbf{A} \mathbf{W}^*, \quad (\text{D.2})$$

with each row of $\mathbf{A} \in \mathbb{R}^{p \times k}$ sampled uniformly on a ball of radius one. We acknowledge choosing initial student weights as linear combinations of the teacher can be artificial and shrinks the first plateau, but our focus on this work was the specialization phase. Nevertheless, this choice and Eq. (D.1) are particularly suitable to theoretical analysis. Once k and p are fixed, the dimension d can be varied without changing \mathbf{Q}^0 , \mathbf{M}^0 and \mathbf{P} , thereby removing any influence of different initial conditions for different d and providing the reader better visualization on the learning curves. To clarify this point, consider the j -th row $\mathbf{w}_j^0 \equiv [\mathbf{W}^0]_j \in \mathbb{R}^d$ of \mathbf{W}^0 :

$$\mathbf{w}_j^0 = \sum_{r=1}^k a_{jr} \mathbf{w}_r^*, \quad (\text{D.3})$$

with $a_{jr} \equiv [\mathbf{A}]_{jr}$. Using Eq. (D.1) one can write

$$q_{jl}^0 = \frac{\mathbf{w}_j^{0\top} \mathbf{w}_l^0}{d} = \sum_{r,r'=1}^k a_{jr} a_{jr'} \underbrace{\frac{\mathbf{w}_r^{*\top} \mathbf{w}_{r'}^*}{d}}_{=\delta_{rr'}} = \sum_{r=1}^k a_{jr} a_{lr}. \quad (\text{D.4})$$

Similarly,

$$m_{jr}^0 = \frac{\mathbf{w}_j^{0\top} \mathbf{w}_r^*}{d} = a_{jr} . \quad (\text{D.5})$$

Thus once \mathbf{A} is fixed, the input dimension d can be varied without affecting the initial conditions. We chose to sample $\mathbf{a}_j \equiv [\mathbf{A}]_j \in \mathbb{R}^k$ on a ball of radius one both to introduce some randomness on the initialization and to keep the initial parameters bounded by one.

We stress that we use these initial conditions to make the data comparable for varying dimension d in the numerical illustrations. Our conclusions do not depend on this particular choice of initial conditions. If one simply takes random initialization $\mathbf{w}_j \sim \mathcal{N}(\mathbf{w}_j | \mathbf{0}, \mathbf{1})$ for each j , the full picture we have presented in this manuscript remains unchanged. In Figure 7 we present an example of curves within the blue region (see Section 3 for the characterization of this regime) with unconstrained Gaussian initialization. Dots represent simulations, while solid lines are obtained by integration of the ODEs given by Eqs. (23), with initial conditions adjusted to match simulations.

Although varying the initial population risk with d slightly changes the exact position where the specialization transition starts, the particular initial conditions adopted in this work do not affect whether the specialization transition takes place or not, comparing to unconstrained Gaussian initialization.

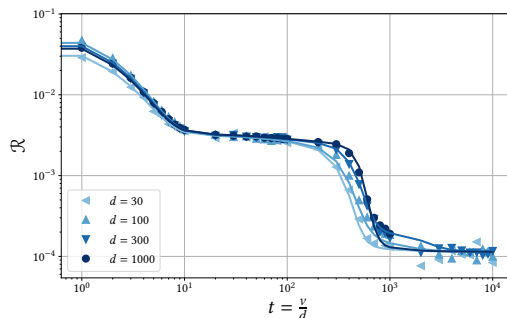


Figure 7: Population risk dynamics for $\kappa = \delta = 0$ (Saad & Solla scaling) : $p_0 = 8, k = 4, \rho_{rs} = \delta_{rs}$. Initialization: $\mathbf{w}_j \sim \mathcal{N}(\mathbf{w}_j | \mathbf{0}, \mathbf{1})$ for $j = 1, \dots, p_0$. Activation function: $\sigma(x) = \text{erf}(x/\sqrt{2})$. Data distribution: $\mathbb{P}(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \mathbf{0}, \mathbf{1})$. Dots represent simulations, while solid lines are obtained by integration of the ODEs given by Eqs. (23), with initial conditions adjusted to match simulations. Observe the difference on the initialization for different d .

References

- [1] S. Mei, A. Montanari, and P.-M. Nguyen, “A mean field view of the landscape of two-layer neural networks,” *Proceedings of the National Academy of Sciences*, vol. 115, no. 33, pp. E7665–E7671, 2018.
- [2] L. Chizat and F. Bach, “On the global convergence of gradient descent for over-parameterized models using optimal transport,” in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc., 2018.
- [3] G. M. Rotskoff and E. Vanden-Eijnden, “Trainability and accuracy of neural networks: An interacting particle system approach,” 2019.
- [4] J. Sirignano and K. Spiliopoulos, “Mean field analysis of neural networks: A central limit theorem,” *Stochastic Processes and their Applications*, vol. 130, no. 3, pp. 1820–1852, 2020.
- [5] D. Saad and S. A. Solla, “On-line learning in soft committee machines,” *Phys. Rev. E*, vol. 52, pp. 4225–4243, Oct 1995.
- [6] S. Goldt, M. Advani, A. M. Saxe, F. Krzakala, and L. Zdeborová, “Dynamics of stochastic gradient descent for two-layer neural networks in the teacher-student setup,” in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32. Curran Associates, Inc., 2019.
- [7] W. Kinzel and P. Ruján, “Improving a network generalization ability by selecting examples,” *Europhysics Letters (EPL)*, vol. 13, no. 5, pp. 473–477, nov 1990.
- [8] O. Kinouchi and N. Caticha, “Optimal generalization in perceptions,” *Journal of Physics A: Mathematical and General*, vol. 25, no. 23, pp. 6243–6250, dec 1992.
- [9] M. Copelli and N. Caticha, “On-line learning in the committee machine,” *Journal of Physics A: Mathematical and General*, vol. 28, no. 6, pp. 1615–1625, mar 1995.
- [10] M. Biehl and H. Schwarze, “Learning by on-line gradient descent,” *Journal of Physics A: Mathematical and General*, vol. 28, no. 3, pp. 643–656, feb 1995.
- [11] P. Riegler and M. Biehl, “On-line backpropagation in two-layered neural networks,” *Journal of Physics A: Mathematical and General*, vol. 28, no. 20, pp. L507–L513, oct 1995.
- [12] D. Saad and S. Solla, “Dynamics of on-line gradient descent learning for multilayer neural networks,” in *Advances in Neural Information Processing Systems*, D. Touretzky, M. C. Mozer, and M. Hasselmo, Eds., vol. 8. MIT Press, 1996.
- [13] R. Vicente, O. Kinouchi, and N. Caticha, “Statistical mechanics of online learning of drifting concepts: A variational approach,” *Machine learning*, vol. 32, no. 2, pp. 179–201, 1998.
- [14] D. Saad, Ed., *On-Line Learning in Neural Networks*, ser. Publications of the Newton Institute. Cambridge: Cambridge University Press, 1999.

- [15] S. Mei, T. Misiakiewicz, and A. Montanari, “Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit,” in *Proceedings of the Thirty-Second Conference on Learning Theory*, ser. Proceedings of Machine Learning Research, A. Beygelzimer and D. Hsu, Eds., vol. 99. PMLR, 25–28 Jun 2019, pp. 2388–2464.
- [16] D. Saad and S. A. Solla, “Exact solution for on-line learning in multilayer neural networks,” *Phys. Rev. Lett.*, vol. 74, pp. 4337–4340, May 1995.
- [17] M. Refinetti, S. Goldt, F. Krzakala, and L. Zdeborova, “Classifying high-dimensional gaussian mixtures: Where kernel methods fail and neural networks succeed,” in *Proceedings of the 38th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, M. Meila and T. Zhang, Eds., vol. 139. PMLR, 18–24 Jul 2021, pp. 8936–8947.
- [18] S. Goldt, B. Loureiro, G. Reeves, F. Krzakala, M. Mézard, and L. Zdeborová, “The gaussian equivalence of generative models for learning with two-layer neural networks,” in *Proceedings of Machine Learning Research*, vol. 145. 2nd Annual Conference on Mathematical and Scientific Machine Learning, 2021, pp. 1–46.
- [19] H. Hu and Y. M. Lu, “Universality laws for high-dimensional learning with random features,” 2020.
- [20] A. Montanari and B. Saeed, “Universality of empirical risk minimization,” 2022.
- [21] C. Wang, Y. C. Eldar, and Y. M. Lu, “Subspace estimation from incomplete observations: A high-dimensional analysis,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 12, no. 6, pp. 1240–1252, 2018.
- [22] Y. Yoshida, R. Karakida, M. Okada, and S.-i. Amari, “Statistical Mechanical Analysis of Online Learning with Weight Normalization in Single Layer Perceptron,” *Journal of the Physical Society of Japan*, vol. 86, no. 4, p. 044002, Apr. 2017.
- [23] P. Del Moral and A. Niclas, “A taylor expansion of the square root matrix function,” *Journal of Mathematical Analysis and Applications*, vol. 465, no. 1, pp. 259–266, 2018.
- [24] B. Aubin, A. Maillard, J. Barbier, F. Krzakala, N. Macris, and L. Zdeborová, “The committee machine: computational to statistical gaps in learning a two-layers neural network,” *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2019, no. 12, p. 124023, dec 2019.
- [25] A. Jacot, F. Gabriel, and C. Hongler, “Neural tangent kernel: Convergence and generalization in neural networks,” in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc., 2018.
- [26] L. Chizat, E. Oyallon, and F. Bach, “On lazy training in differentiable programming,” in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32. Curran Associates, Inc., 2019.
- [27] F. Bach and L. Chizat, “Gradient descent on infinitely wide neural networks: Global convergence and generalization,” *arXiv preprint arXiv:2110.08084*, 2021.
- [28] Y. S. Tan and R. Vershynin, “Phase retrieval via randomized Kaczmarz: theoretical guarantees,” *Information and Inference: A Journal of the IMA*, vol. 8, no. 1, pp. 97–123, 04 2018.

- [29] G. B. Arous, R. Gheissari, and A. Jagannath, “Online stochastic gradient descent on non-convex losses from high-dimensional inference,” *Journal of Machine Learning Research*, vol. 22, no. 106, pp. 1–51, 2021.
- [30] —, “Algorithmic thresholds for tensor PCA,” *The Annals of Probability*, vol. 48, no. 4, pp. 2052 – 2087, 2020.
- [31] C. Wang, H. Hu, and Y. Lu, “A solvable high-dimensional model of gan,” in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32. Curran Associates, Inc., 2019.