

## Lecture 12: Mean Field Limit

### 1 Finishing up the NTK Analysis

#### 1.1 Recalling the Setting

We are given a dataset  $(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$  and a student network  $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ . We are trying to learn parameters  $\theta \in \mathbb{R}^p$ , where  $\theta$  is initialized to some  $\theta_0$ . We define the loss of some function  $g$  to be  $\hat{L}(g) \triangleq \frac{1}{2} \|g(X) - y\|_2^2$  and  $\hat{L}_0 \triangleq \hat{L}(\gamma f_{\theta_0})$ , where  $\gamma > 0$  is a scaling parameter.

We can define a **gradient flow** on the parameter space  $\Theta$  through the equation

$$d\theta_t = -\nabla_\theta \hat{L}(\gamma f_{\theta_t}) dt = -\gamma J_t^T \nabla \hat{L}(\gamma f_{\theta_t}) dt,$$

where the **Jacobian** matrix is given by

$$J_t = J_{\theta_t} = \begin{pmatrix} \nabla_\theta f_{\theta_t}(x_1) \\ \vdots \\ \nabla_\theta f_{\theta_t}(x_n) \end{pmatrix} \in \mathbb{R}^{n \times p}$$

The main idea is to compare the gradient flow to **linearized dynamics**, which is given by the equations

$$\begin{aligned} f_\theta^{lin}(x) &= f_{\theta_0}(x) + J_0(\theta - \theta_0) \\ d\tilde{\theta}_t &= -\nabla_\theta \hat{L}(\gamma f_{\tilde{\theta}_t}^{lin}) = -\gamma J_0^T \nabla \hat{L}(\gamma f_{\tilde{\theta}_t}^{lin}), \end{aligned}$$

where we note here that the Jacobian does not change in time for the linearized network. Two assumptions from last time were

1.  $J_\theta$  is Lipschitz in  $\theta$  with some constant  $\beta$ .
2.  $J_0 = J_{\theta_0}$  is full rank.

Before we complete the analysis, we recall the lemmas we proved last time.

**Lemma 1.** Suppose that  $Q(t) \geq \lambda \cdot Id$  for all  $t$ . Then for  $(g_t)$  given by

$$dg_t = -Q(t) \nabla \hat{L}(g_t) dt,$$

we have

$$\hat{L}(g_t) \leq \hat{L}(g_0) \cdot \exp(-2\lambda t).$$

In lecture, the way this was written was

$$\hat{L}(\tilde{\theta}_t) \leq \exp(-\Omega(\sigma_{\min}^2(J_0)\gamma^2 t)) \hat{L}_0$$

*Proof.* The proof idea was through a strong convexity style argument. □

**Lemma 2 (Lemma 2).** If  $\theta$  is close to  $\theta_0$ , then  $J_\theta$  is close to  $J_{\theta_0}$ .

*Proof.* The idea of the proof was the triangle inequality and the assumption that  $J$  is Lipschitz, that is,

$$\|J_\theta - J_{\theta_0}\|_{op} \leq \beta \|\theta - \theta_0\|_2$$

for some  $\beta > 0$ . □

Today, we will do the final step, proving

**Lemma 3.** *Suppose that  $(\hat{\theta}_t)$  satisfies*

$$d\hat{\theta}_t = -S(t)^T \nabla \hat{L}(g_{\hat{\theta}_t})$$

for some network  $G$ , and that  $\underline{\lambda} \cdot Id \leq S(t)S(t)^T \leq \bar{\lambda} \cdot Id$ . Then

$$\|\hat{\theta}_t - \hat{\theta}_0\|_2 \leq \frac{\sqrt{\bar{\lambda}}}{\underline{\lambda}} \|g_{\hat{\theta}_0}(x) - y\|.$$

The key feature here is that the numerator features a square root.

*Proof.*

$$\begin{aligned} \hat{\theta}_t &= \hat{\theta}_0 - \int_0^t S(s)^T \nabla \hat{L}(g_{\hat{\theta}_s}) ds \\ \Rightarrow \|\hat{\theta}_t - \hat{\theta}_0\| &= \left\| \int_0^t S(s)^T \nabla \hat{L}(g_{\hat{\theta}_s}) ds \right\| \\ &\leq \int_0^t \underbrace{\|S(s)\|_{op}}_{\leq \sqrt{\bar{\lambda}}} \cdot \underbrace{\|\nabla \hat{L}(g_{\hat{\theta}_s})\|}_{g_{\hat{\theta}_s}(x) - y} ds \\ &\leq \sqrt{\bar{\lambda}} \cdot \int_0^t \underbrace{\|g_{\hat{\theta}_s}(x) - y\|}_{\leq \exp(-\lambda s) \cdot \|g_{\hat{\theta}_0} - y\| \text{ (by previous lemma)}} ds \\ &\leq \sqrt{\bar{\lambda}} \cdot \|g_{\hat{\theta}_0} - y\| \cdot \int_0^t \exp(-\lambda s) ds \\ &\leq \sqrt{\bar{\lambda}} \cdot \|g_{\hat{\theta}_0} - y\| \cdot \frac{1}{\lambda} \end{aligned}$$

Applying this inequality to  $\hat{\theta}_t = \tilde{\theta}_t$ ,  $g_\theta = f_\theta^{lin}$ ,  $S(t) = \gamma J_0$ , we have that

$$\begin{aligned} \|\tilde{\theta}_t - \theta_0\| &\leq \frac{\sqrt{\sigma_{\max}^2(J_0)}}{\sigma_{\min}^2(J_0)} \cdot \|f_{\theta_0}(x) - y\| \\ &\leq \frac{\sqrt{2\gamma^2 \sigma_{\max}^2(J_0)}}{\gamma^2 \sigma_{\min}^2(J_0)} \cdot \sqrt{\hat{L}_0} \\ &\lesssim \frac{\sigma_{\max}(J_0)}{\gamma \sigma_{\min}^2} \cdot \sqrt{\hat{L}_0} \end{aligned}$$

Then, substituting our bounds on the eigenvalues of  $J_0$  suffices for the proof. □

## 2 Mean-Field Limit

To recall, in the NTK regime, we have as the trainer network

$$f_\theta(x) = \gamma \sum_{i=1}^N a_i \sigma(\langle w_i, x \rangle),$$

and the NTK regime is characterized as  $\gamma \gg 1/N$ . As showed previously, the Jacobian does not change much over this scaling, and so the network is well-approximated by its Taylor expansion around the parameters at

initialization. However, gradient descent in the NTK regime is bottlenecked by what kernel methods can do. Then, the natural question is what happens in the regime  $\gamma \asymp \frac{1}{N}$ , which is known as the *mean field regime*.

We first introduce notation, then explain the defining features of the mean field limit. We first introduce the student network

$$f_\theta(x) = \frac{1}{N} \sum_{i=1}^N a_i \sigma(\langle w_i, x \rangle)$$

It will be convenient to denote this by  $f_\theta(x) = \frac{1}{N} \sum_{i=1}^N \sigma(x; \theta_i)$  We will consider population gradient descent:

$$\theta^{(k+1)} \leftarrow \theta^{(k)} - \eta_k \nabla L(\theta^{(k)})$$

where  $L(\theta) = \mathbb{E}_{x,y} [(y - f_\theta(x))^2]$  is the test loss (can be approximated using online gradient descent). This loss function  $L(\theta)$  can be decomposed as

$$L(\theta) = \mathbb{E}[y^2] + \frac{2}{N} \sum_{i=1}^N V(\theta_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\theta_i, \theta_j),$$

where

$$V(\theta_i) \triangleq -\mathbb{E}[y \cdot \sigma(x; \theta_i)]$$

is known as the external field and

$$U(\theta_i, \theta_j) \triangleq \mathbb{E}[\sigma(x; \theta_i) \cdot \sigma(x; \theta_j)]$$

are "pairwise interactions" between particles/neurons. Some convenient notation will be to set

$$\nabla_{\theta_i} L(\theta) = \frac{2}{N} \nabla \Psi_\theta(\theta_i) \quad \text{for} \quad \Psi_\theta(\theta_i) \triangleq V(\theta_i) + \frac{1}{N} \sum_{j=1}^N U(\theta_i, \theta_j).$$

## 2.1 Physics Intuition for the Mean Field Limit

The physics intuition is to regard each  $\theta_i$  is an interacting particle. The idea is that as  $N \rightarrow \infty$ , the fluctuations in the "environment" around any given particle average out, so every particle experiences same "average environment," hence the name mean field. In the mean field limit, at any point in time, all particles are i.i.d. draws from same distribution. The natural question, hence, is how does this distribution evolve over the course of training?

Let's consider the loss

$$L(\theta) = \mathbb{E}[y^2] + \frac{2}{N} \sum_{i=1}^N V(\theta_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\theta_i, \theta_j)$$

We can think of the term  $\frac{2}{N} \sum_{i=1}^N V(\theta_i)$  as similar-valued to  $2\mathbb{E}[V(\theta_i)]$  for  $\theta_i$  drawn from the "empirical distribution"  $\hat{\rho}_\theta \triangleq \frac{1}{N} \sum_i \delta_{\theta_i}$ , and the second term  $\frac{1}{N^2} \sum_{i,j=1}^N U(\theta_i, \theta_j)$  can similarly be thought of as some  $U$ -Statistic.

Hence, the idea is that instead of parametrizing in terms of  $\theta = (\theta_1, \dots, \theta_N)$ , we can parametrize in terms of a probability distribution  $\rho$  over  $\mathbb{R}^d$ . Then  $L$  becomes

$$L(\rho) \triangleq \mathbb{E}[y^2] + 2 \int V(\theta) d\rho(\theta) + \int U(\theta, \theta') d\rho(\theta) d\rho(\theta')$$

To give some more intuition, the basic idea in statistical physics is that when the number of particles in a system is very large, instead of considering the dynamics of all of the individual particles, we can instead consider some probability distribution induced by the particles. This assumption is justified in the  $N \rightarrow \infty$  limit due to a law of large numbers argument.

Let's do a comparison of what the dynamics look like in the finite  $N$  case and in the mean-field limit. In the finite  $N$  case, the loss function,  $\Psi$  (parameterized by  $\theta$ ), and the dynamics are given by

$$\begin{aligned} L(\theta) &= \mathbb{E} [y^2] + \frac{2}{N} \sum_{i=1}^N V(\theta_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\theta_i, \theta_j) \\ \Psi_\theta(\theta_i) &\triangleq V(\theta_i) + \frac{1}{N} \sum_{j=1}^N U(\theta_i, \theta_j) \\ d\theta_i^t &= -\nabla \Psi_\theta(\theta_i^t) dt, \quad \theta_i^0 \sim \rho_0. \end{aligned}$$

In the mean-field limit, these equations are given by

$$\begin{aligned} L(\rho) &\triangleq \mathbb{E} [y^2] + 2 \int V(\theta) d\rho(\theta) + \int U(\theta, \theta') d\rho(\theta) d\rho(\theta') \\ \Psi_\rho(\theta) &\triangleq V(\theta) + \int U(\theta, \theta') d\rho(\theta') \\ d\bar{\theta}_i^t &= -\nabla \Psi_{\rho_t}(\bar{\theta}_i^t) dt, \quad \rho_t = \text{law}(\bar{\theta}_i^t) \end{aligned}$$

One difference here is that the parameters in question are not individual  $\theta$ 's, but rather probability distributions. Secondly, the function  $\Psi$  in the dynamics is now parameterized by  $\rho_t$ . The idea here is that the dynamics of the mean-field are now parameterized by the mean-field variables themselves. In particular, one can show that the dynamics are given by a continuity PDE:

$$\partial_t \rho_t = \text{Div}(\rho_t \cdot \nabla \Psi_{\rho_t}),$$

which holds in the weak sense. In the above,  $-\nabla \Psi_{\rho_t}$  is called the velocity field of  $\rho_t$ , and the idea is that the process  $\rho_t$  is performing gradient descent in the space of probability distributions (equipped with the Wasserstein metric) with respect to the function  $L(\rho)$ . There is a rich theory about this material; one good source is [AGS06].

## 2.2 Derivation of the Continuity Equation

Recall again the continuity equation given above:

$$\partial_t \rho_t = \text{Div}(\rho_t \cdot \nabla \Psi_{\rho_t}).$$

In this statement, we mean that the PDE holds in the weak sense, in that for any "nice" (e.g. bounded, differentiable, with bounded gradient) test function  $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}$ , we have

$$\int \varphi(\theta) \partial_t \rho_t(\theta) d\theta = \int \varphi(\theta) \cdot \text{div}(\rho_t \cdot \nabla \Psi_{\rho_t})(\theta) d\theta \quad (\star)$$

We formulate this as being satisfied in the weak sense because differentiable solutions to  $(\star)$  might not exist. We now show the derivation for this equation. Note that for  $\bar{\theta}_t \sim \rho_t$ ,

$$\begin{aligned} \text{LHS of } (\star) &= \frac{\partial}{\partial t} \mathbb{E} \left[ \varphi(\bar{\theta}_t) \right] \\ &= \mathbb{E} \left[ \left\langle \nabla \varphi(\bar{\theta}_t), \frac{d}{dt} \bar{\theta}_t \right\rangle \right] \quad (\text{differentiate under integral}) \\ &= \int \langle \nabla \varphi(\theta), -\nabla \Psi_{\rho_t}(\theta) \rangle d\rho_t(\theta) \quad (\text{gradient flow for } \bar{\theta}_t) \\ &= - \int \langle \nabla \varphi(\theta), \nabla \Psi_{\rho_t}(\theta) \rangle d\rho_t(\theta) \\ &= \text{RHS of } (\star) \quad (\text{integration by parts}) \end{aligned}$$

### 2.3 Non-asymptotic Convergence to the Mean-Field Limit

We use a method known as the "Propagation of chaos." Some original references are by Kac in 1956, McKean in 1969, and Sznitman in 1991; a comprehensive reference is in [CD22]. We want to compare  $(\theta_i^{(k)})_{k=0,1,2,\dots}$  and  $(\bar{\theta}_i^t)_{t \geq 0}$ , where the first are the gradient descent iterates given by  $\theta_i^{(k+1)} \leftarrow \theta_i^{(k)} - h \nabla L(\theta^{(k)})$ , and the second are the mean-field iterates given by  $d\bar{\theta}_i^t = -\nabla L_{\rho_t}(\bar{\theta}_i^t) dt$ , where  $\rho_t = \text{law}(\bar{\theta}_i^t)$ . Note that

$$\begin{aligned}\theta_i^{(k)} &= \theta_i^{(0)} + 2h \sum_{l=0}^{k-1} F_i(\theta^{(l)}; (x_{l+1}, y_{l+1})) \\ \bar{\theta}_i^t &= \bar{\theta}_i^{(0)} + 2 \int_0^t G(\bar{\theta}_i^s; \rho_s) ds\end{aligned}$$

for  $F_i(\theta; (x, y)) \triangleq (y - f_\theta(x)) \cdot \nabla_{\theta_i} \sigma(x; \theta_i)$ ,  $G(\theta; \rho) \triangleq -\nabla \Psi_\rho(\theta)$ . The ultimate goal is to upper bound  $\|\bar{\Theta}_i^{kh} - \theta_i^{(k)}\|$ . To do so, will bound by a self-similar expression of the form

$$(\text{small terms}) + \int_0^{kh} \|\bar{\theta}_i^s - \theta_i^{(\lfloor s/h \rfloor)}\| ds.$$

This will imply (by Grönwall's inequality), the desired bound.

Now, let  $[S] = h \cdot \lfloor s/h \rfloor$  and consider

$$\begin{aligned}\|\bar{\theta}_i^{kh} - \theta_i^k\| &= 2 \left\| \int_0^{kh} G(\bar{\theta}_i^s; \rho_s) ds - h \sum_{l=0}^{k-1} F_i(\theta^{(l)}; (x_{l+1}, y_{l+1})) \right\| \\ &\leq 2 \left\| \int_0^{kh} \left[ G(\bar{\theta}_i^s; \rho_s) - G(\bar{\theta}_i^{[s]}; \rho_{[s]}) \right] ds \right\| \quad (1) \\ &\quad + 2 \left\| \int_0^{kh} \left[ G(\bar{\theta}_i^{[s]}; \rho_{[s]}) - G(\theta_i^{(\lfloor s/h \rfloor)}; \rho_{[s]}) \right] ds \right\| \quad (2) \\ &\quad + 2 \left\| h \sum_{l=0}^{k-1} \left[ G(\theta_i^{(l)}; \rho_{lh}) - F_i(\theta^{(l)}; (x_{l+1}, y_{l+1})) \right] \right\| \quad (3)\end{aligned}$$

We first bound term (1) (easy). This term is small because because  $G$  is Lipschitz by assumption, and we can show  $f$  varies smoothly over time so that  $\rho_s$  and  $\rho_{[s]}$  are close. Term (2) is bounded by the Lipschitzness of  $G$ :

$$\left\| G(\bar{\theta}_i^s; \rho_{[s]}) - G(\theta_i^{(\lfloor s/h \rfloor)}; \rho_{[s]}) \right\| \lesssim \|\bar{\theta}_i^s - \theta_i^{(\lfloor s/h \rfloor)}\|,$$

so (2) is bounded by

$$\int_0^{kh} \underbrace{\|\bar{\theta}_i^s - \theta_i^{(\lfloor s/h \rfloor)}\|}_{\substack{\text{looks analogous} \\ \text{to what we want} \\ \text{to bound on the LHS}}} ds$$

Now to bound (3), consider

$$\sum_{l=0}^{k-1} \left[ G(\theta_i^{(l)}; \rho_{lh}) - F_i(\theta^{(l)}; (x_{l+1}, y_{l+1})) \right].$$

The key idea is to note that  $F_i(\theta^{(l)}; (x_{l+1}, y_{l+1}))$  has expectation  $G(\theta_i^{(l)}; \hat{\rho}_l)$ , where  $\hat{\rho}_l$  is the empirical distribution of  $\frac{1}{N} \sum_{i=1}^N \delta_{\theta_i^{(l)}}$ . Then, over many steps  $l$ , the total deviation between the  $F_i(\theta^{(l)}; (x_{l+1}, y_{l+1}))$  and the  $G(\theta_i^{(l)}; \hat{\rho}_l)$  is of order  $h\sqrt{k\rho}$  by Martingale concentration. Then, replacing  $F_i$  with its expectation, it remains to bound

$$\begin{aligned}
& \sum_{l=0}^{k-1} \left[ G\left(\theta_i^{(l)}; \rho_{lh}\right) - G\left(\theta_i^{(l)}; \hat{\rho}_l\right) \right] \\
&= \frac{1}{N} \sum_{l=0}^{k-1} \sum_{j=1}^N \left[ \mathbb{E}_{\bar{\theta}} U\left(\theta_i^{(l)}, \bar{\theta}_j^{lh}\right) - U\left(\theta_i^{(l)}, \theta_j^{(l)}\right) \right]
\end{aligned}$$

Again, by Martingale concentration we can essentially replace  $\mathbb{E}_{\bar{\theta}} U\left(\theta_i^{(l)}, \bar{\theta}_j^{lh}\right)$  (deterministic) with  $U\left(\theta_i^{(l)}, \bar{\theta}_j^{lh}\right)$  (random). Then, we use Lipschitzness of  $U$  to get

$$\begin{aligned}
& \frac{1}{N} \sum_{l=0}^{k-1} \sum_{j=1}^N \left\| U\left(\theta_i^{(l)}, \bar{\theta}_j^{lh}\right) - U\left(\theta_i^{(l)}, \theta_j^{(l)}\right) \right\| \\
& \leq \frac{1}{N} \sum_{l=0}^{k-1} \sum_{j=1}^N \left\| \bar{\theta}_j^{lh} - \theta_j^{(l)} \right\|
\end{aligned}$$

Once again, this last term looks similar to what we want to bound. This yields the self-similar equation we sought after.

## 2.4 Note on the PDE when the Data Distribution has Symmetries

When data distribution has symmetries, PDE simplifies considerably. Suppose that the training data  $\{(x_i, y_i)\}$  satisfies  $x_i \sim N(0, I)$  and  $y_i = \varphi(\Pi x)$  where  $\Pi$  is a projection to a low-dimensional subspace  $V^*$ .

Then the joint dist over  $(x, y)$  is invariant under rotations of  $x$  that preserve  $V^*$ , ie.  $Rv \in V^* \forall v \in V^*$ . Consider this observation: Let  $R$  be such a rotation. If  $\rho_0$  and  $\rho'_0$  are two different initializations of the weights related by  $\rho'_0 = R_{\#} \rho_0$  (i.e., to sample  $(a', w')$  from  $\rho'_0$ , one can sample  $(a, w)$  from  $\rho_0$  and then take  $a' = a, w' = R w$ ), then  $\rho'_t = R_{\#} \rho_t$ .

Hence, if  $\rho_0$  is rotation-invariant, then  $\rho_t$  is invariant to rotations preserving  $V^*$  for any  $t \geq 0$ !  $\rho_t$  is thus completely specified by the distribution on

$$\left( a, \underbrace{\Pi w}_{\vec{s}}, \underbrace{\|\Pi^1 w\|_2}_r \right),$$

i.e., we get a  $\dim(V^*) + 2$  dimensional PDE! Then denoting the distribution, or  $(a, \vec{s}, r)$  by  $\bar{\rho}_t$ , we have

$$\partial_t \bar{\rho}_t = \operatorname{div}(\bar{\rho}_t \cdot \nabla_{\vec{s}} \Psi_{\bar{\rho}_t}) + \partial_a(\bar{\rho}_t \cdot \partial_a \Psi_{\bar{\rho}_t}) + \frac{1}{r} \partial_r(r \cdot \bar{\rho}_t \cdot \partial_r \bar{\Psi}_{\bar{\rho}_t}).$$

## 3 Analysis of the Mean-Field Theory

Recall that the game plan for understanding training dynamics of overparameterized neural networks was

1. Define the limiting object.
2. Show that we quickly converge to the limiting object as  $N \rightarrow \infty$ .
3. Prove optimization/generalization guarantees for limiting object.

### 3.1 Theorems Regarding Speed of Convergence

Works such as [MMN18] and [MMM19] address this second point. One theorem that they proved is:

**Theorem 1** (Mei-Montanari-Nguyen '18). *Assumptions: (1)  $\nabla V, \nabla U$  bounded Lipschitz, (2)  $\sigma$  bounded, (3)  $\nabla_{\theta} \sigma(x; \theta)$  has sub-Gaussian tails.*

*Then let  $(\theta^{(k)})_{k=0,1,2,\dots}$  denote iterates of gradient descent with step size  $h$  and let  $(\bar{\theta}^t)_{t \geq 0}$  denote mean-field gradient flow. Then with probability  $1 - \delta$  in the randomness of the initialization and the training examples,*

$$\begin{aligned} & \max_{i \in [N]} \sup_{k=0, \dots, T/h} \left\| \theta_i^{(k)} - \bar{\theta}_i^{kh} \right\|_2 \\ & \lesssim e^{O(T)} \cdot \sqrt{\max(1/N, h)} \cdot [\sqrt{p + \log(N \max(1, T/h))} + \sqrt{\log 1/\delta}] \end{aligned}$$

### 3.2 Theorems Regarding Asymptotics

Regarding the third point above, it is in general quite difficult to prove things about the PDE. The current understanding is largely limited to asymptotics, toy examples (e.g., generalized linear models), and experiments. One example is from [CB18], from which there is the following information theorem about the evolution of  $L(\rho_t)$  in noiseless gradient descent.

**Theorem 2.** *Suppose that  $\sigma$  is sigmoid or ReLU, and the distribution over  $x$  has finite fourth-order moments. If the support of the random initialization is chosen appropriately and the distribution over  $x$  satisfies certain a certain ‘‘Sard-type regularity’’ condition, then if  $\rho_t \rightarrow \rho_{\infty}$ , then  $\rho_{\infty}$  is a global minimizer of  $L(\rho)$ .*

The proof of this theorem exploits facts from the Wasserstein gradient flow, as well as homogeneity/partial homogeneity of the activation function. Follow-up works include [NP23] and [EW21].

Another theorem regarding asymptotics is the evolution of  $L(\rho_t)$ .

**Theorem 3** ([MMN18] (informal)). *The theorem extends the mean-field picture to ‘‘noisy gradient descent’’, i.e.  $\theta \leftarrow (1 - 2\lambda\eta) \cdot \theta - 2\eta \cdot \nabla L(\theta) + N(0, \eta/\beta)$ . As  $t \rightarrow \infty$ , the resulting continuity PDE converges to minimizer of the regularized loss (free energy)*

$$F_{\beta, \lambda}(\rho) \stackrel{\text{def}}{=} L(\rho)/2 + (\lambda/2) \cdot \mathbb{E}_{\rho} [\|\theta\|^2] - \beta^{-1} \text{Ent}(\rho)$$

*In fact, limit distribution satisfies*

$$\rho_{\infty}(\theta) \propto \exp\left(-\beta \Psi_{\rho_{\infty}}^{(\lambda)}(\theta)\right)$$

### 3.3 Toy Models

One setting in which the continuity PDE greatly simplifies is when the underlying data distribution has symmetries. For instance, suppose that the  $x$ 's are Gaussian and that  $y = \phi(\langle w^*, x \rangle)$  (the single index model), so that the function secretly only depends on a 1D subspace. Recall this is a setting that kernel methods perform terribly at because they fail to efficiently learn the relevant feature  $\langle w^*, \cdot \rangle$ . Then, because data distribution is invariant to any rotation that preserves  $w^*$ , the PDE simplifies dramatically! In such settings where symmetries drastically reduce the dimension of the PDE, we can numerically solve it and obtain sharp predictions. The following is an example from [ABAM22]. Another example is from [MMN18]. In this paper, the context is classifying isotropic Gaussians. The data is a mixture of

$$\begin{cases} x \sim N(0, (1 + \Delta)^2 \cdot \text{Id}_d) & y = 1 \\ x \sim N(0, (1 - \Delta)^2 \cdot \text{Id}_d) & y = -1 \end{cases}$$

In this case, the reduced PDE is 1-dimensional (only need to track distribution of  $\|\theta\|_2$ ), and they obtain rigorous end-to-end guarantees for this problem.

Another example is [BMZ23]. In this paper, the continuity PDE is for learning  $f(x) = \phi(\langle w^*, x \rangle)$  for  $\phi(z) = \text{He}_0(z) - \text{He}_1(z) + \frac{2}{3}\text{He}_2(z)$  using a one-hidden-layer ReLU network over Gaussian examples. They demonstrate two empirically observed phenomena: (1) plateaus in the loss curve, interspersed with sharp drops and (2) longer and longer time scales. A key (partially rigorous) finding is that the mean field gradient flow incrementally learns low-degree Hermite components of single index models.

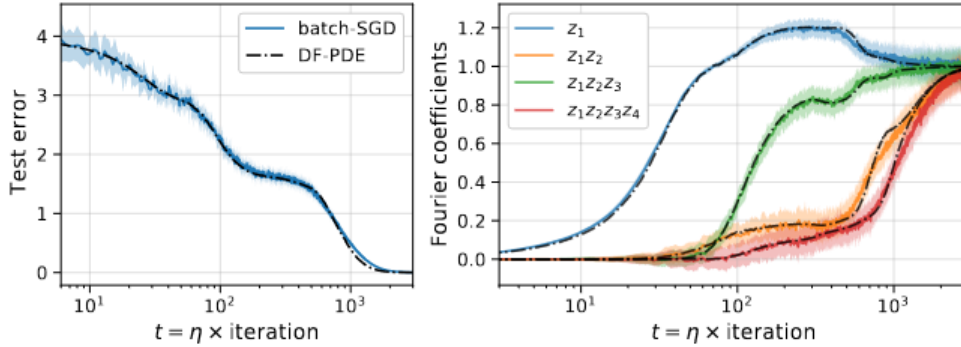


Figure 1: Comparison between (bSGD) and (DF-PDE) dynamics for  $h_*(z) = z_1 + z_1 z_2 + z_1 z_2 z_3 + z_1 z_2 z_3 z_4$ . Left: Test error. Right: Fourier coefficients of  $\hat{f}_{\text{NN}}(x; \Theta^{t/\eta})$  and  $\hat{f}_{\text{NN}}(z; \bar{\rho}_t)$ . The dashed-dotted black lines correspond to (DF-PDE) and the continuous colored line to (bSGD). The test errors and Fourier coefficients are evaluated with  $m = 300$  test samples and for (bSGD), we report the average and 95% confidence interval over 10 experiments.

Currently, the only “end-to-end” non-asymptotic, global convergence result for learning single-index models with “standard” gradient descent (no layerwise training, no weird learning rate schedules, no sharp gradient clipping, etc.) is

**Theorem 4** ([MHD<sup>+</sup>23] (informal)). *Projected gradient descent on a one-hidden-layer student network with quartic polynomial activations and polynomial width learns certain functions of the form  $y = p(\langle w^*, x \rangle)$ , where  $p$  is a degree-4 polynomial, over Gaussian inputs using  $O(d^{3.1})$  samples.*

### 3.4 CSQ Revisited

It’s helpful to keep in mind what is generally possible (by any algorithm). Consider single index models ( $y = \phi(\langle w^*, x \rangle)$ ), where the complexity of CSQ algorithms is dictated by the “information exponent,” which is the smallest  $s$  for which the  $s$ -th Hermite coefficient of  $\phi$  is nonzero.

As examples, [AGJ21] shows that online SGD on single neuron learns in time  $d^S$ . [ABAM23] generalizes (“leap complexity”) to multi-index models ( $y = \phi(\Pi_W x)$ ), layerwise training of overparametrized model learns in time  $d^{\text{leap}}$ . These are optimal for CSQ algorithms (next unit: lower bound), but there are more efficient non-CSQ algorithms (filtered PCA) that achieve  $O(d)$  sample complexity and (fixed) polynomial runtime [CM20].

## References

- [ABAM22] Emmanuel Abbe, Enric Boix-Adsera, and Theodor Misiakiewicz. The merged-staircase property: a necessary and nearly sufficient condition for sgd learning of sparse functions on two-layer neural networks, 2022.
- [ABAM23] Emmanuel Abbe, Enric Boix-Adsera, and Theodor Misiakiewicz. Sgd learning on neural networks: leap complexity and saddle-to-saddle dynamics, 2023.
- [AGJ21] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient descent on non-convex losses from high-dimensional inference, 2021.
- [AGS06] L. Ambrosio, N. Gigli, and G. Savare. *Gradient Flows: In Metric Spaces and in the Space of Probability Measures*. Lectures in Mathematics. ETH Zürich. Birkhäuser Basel, 2006.
- [BMZ23] Raphaël Berthier, Andrea Montanari, and Kangjie Zhou. Learning time-scales in two-layers neural networks, 2023.



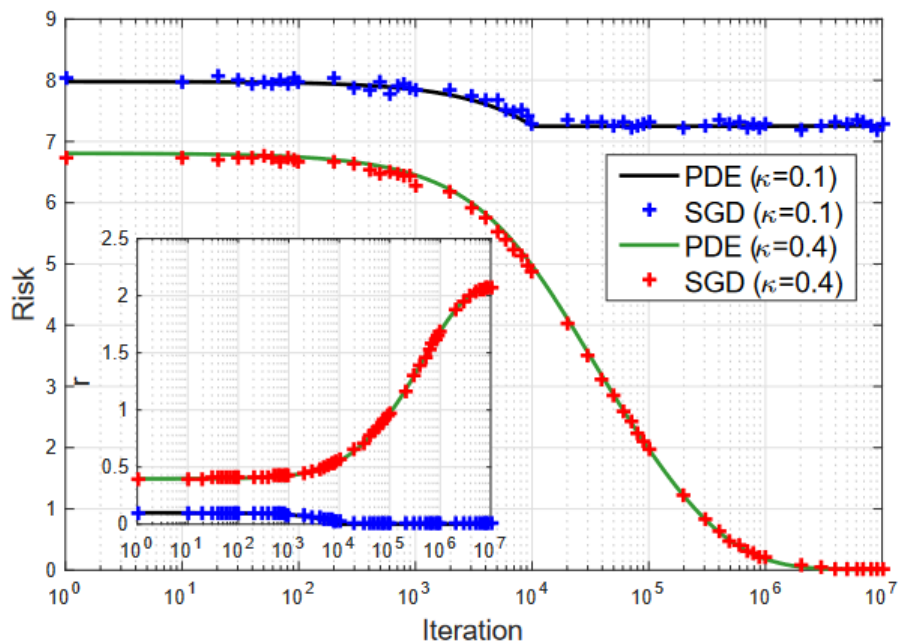


Figure 4: Separating two isotropic Gaussians, with a non-monotone activation function (see text for details). Here  $N = 800$ ,  $d = 320$ ,  $\Delta = 0.5$ . The main frame presents the evolution of the population risk along the SGD trajectory, starting from two different initializations of  $(w_i^0)_{i \leq N} \sim_{iid} \mathbf{N}(0, \kappa^2/d \cdot \mathbf{I}_d)$  for either  $\kappa = 0.1$  or  $\kappa = 0.4$ . In the inset, we plot the evolution of the average of  $\|w\|_2$  for the same conditions. Symbols are empirical results. Continuous lines are prediction obtained with the reduced PDE (13).

- [CB18] Lenaic Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport, 2018.
- [CD22] Louis-Pierre Chaintron and Antoine Diez. Propagation of chaos: A review of models, methods and applications. *Kinetic and Related Models*, 15(6):1017, 2022.
- [CM20] Sitan Chen and Raghu Meka. Learning polynomials of few relevant dimensions, 2020.
- [EW21] Weinan E and Stephan Wojtowytsch. Representation formulas and pointwise properties for barron functions, 2021.
- [MHD<sup>+</sup>23] Arvind Mahankali, Jeff Z. Haochen, Kefan Dong, Margalit Glasgow, and Tengyu Ma. Beyond ntk with vanilla gradient descent: A mean-field analysis of neural networks with polynomial width, samples, and time, 2023.
- [MMM19] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In Alina Beygelzimer and Daniel Hsu, editors, *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 2388–2464. PMLR, 25–28 Jun 2019.
- [MMN18] Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33), jul 2018.
- [NP23] Phan-Minh Nguyen and Huy Tuan Pham. A rigorous framework for the mean field limit of multilayer neural networks, 2023.

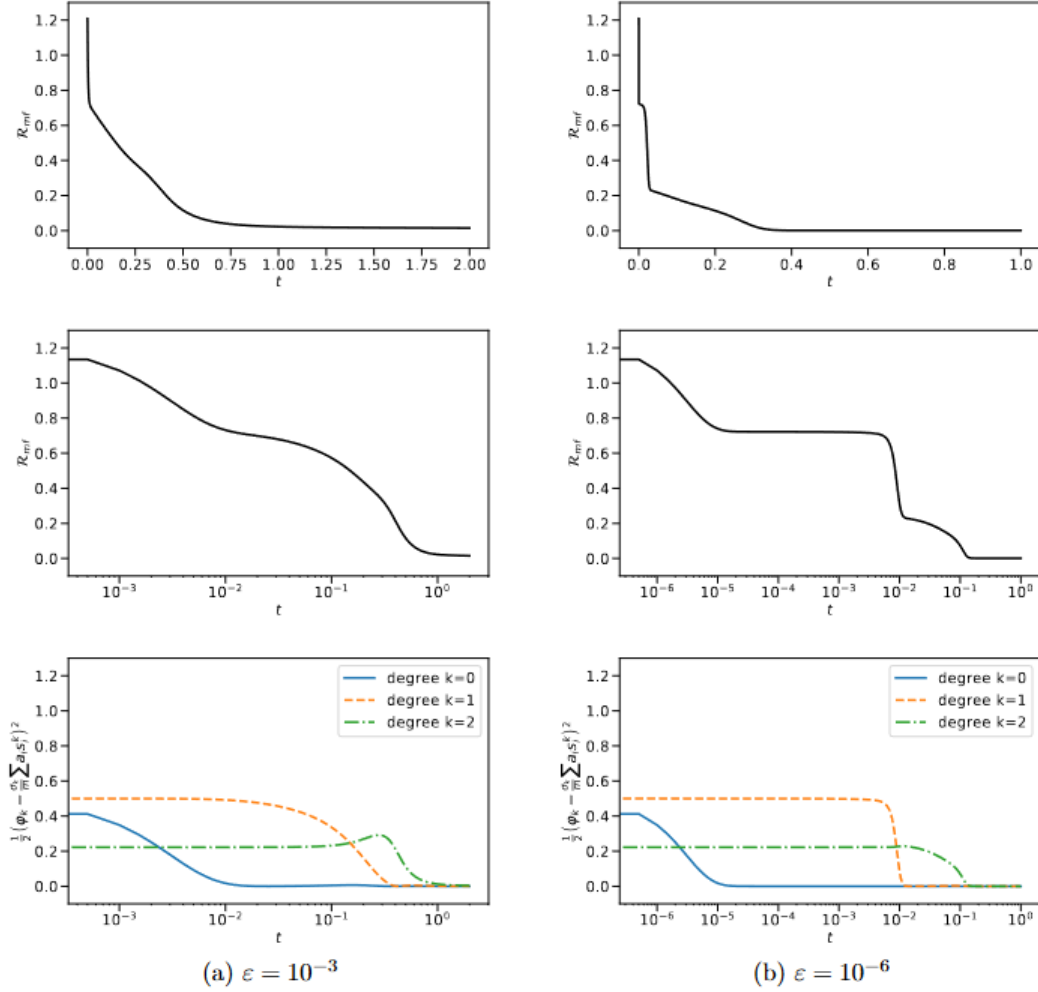


Figure 2: Simulation of the mean field neuron dynamics of Eqs. (23), with the target function of Eq. (37) and ReLU activations. We use learning rate ratios  $\varepsilon = 10^{-3}$  (left) and  $\varepsilon = 10^{-6}$  (right) and we use  $m = 10$  neurons. First two rows: evolution of the risk  $\mathcal{R}_{mf}$  of Eq. (22), in linear and log-scales. Third row: evolution of the first three terms of the sum of (38).