

Chapter 11

Wide neural networks

In this chapter we explore the dynamics of training neural networks of large width. Throughout we focus on the situation where we have data pairs

$$(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R} \quad i \in \{1, \dots, m\}, \quad (11.0.1a)$$

and wish to train a neural network $\Phi(\mathbf{x}, \mathbf{w})$ depending on the input $\mathbf{x} \in \mathbb{R}^d$ and the parameters $\mathbf{w} \in \mathbb{R}^n$, by minimizing the square loss objective defined as

$$f(\mathbf{w}) := \sum_{i=1}^m (\Phi(\mathbf{x}_i, \mathbf{w}) - y_i)^2, \quad (11.0.1b)$$

which is a multiple of the empirical risk $\widehat{\mathcal{R}}_S(\Phi)$ in (1.2.3) for the sample $S = (\mathbf{x}_i, y_i)_{i=1}^m$ and the square-loss. We exclusively focus on gradient descent with a constant step size h , which yields a sequence of parameters $(\mathbf{w}_k)_{k \in \mathbb{N}}$. We aim to understand the evolution of $\Phi(\mathbf{x}, \mathbf{w}_k)$ as k progresses. For linear mappings $\mathbf{w} \mapsto \Phi(\mathbf{x}, \mathbf{w})$, the objective function (11.0.1b) is convex. As established in the previous chapter, gradient descent then finds a global minimizer. For typical neural network architectures, $\mathbf{w} \mapsto \Phi(\mathbf{x}, \mathbf{w})$ is not linear, and such a statement is in general not true.

Recent research has highlighted that neural network behavior tends to linearize in the parameters as network width increases [105]. This allows to transfer some of the results and techniques from the linear case to the training of neural networks. We start this chapter in Sections 11.1 and 11.2 by recalling (kernel) least-squares methods, which describe linear (in \mathbf{w}) models. Following [129], the subsequent sections explore why in the infinite width limit neural networks exhibit linear-like behavior. In Section 11.5.2 we formally introduce the linearization of $\mathbf{w} \mapsto \Phi(\mathbf{x}, \mathbf{w})$. Section 11.4 presents an abstract result showing convergence of gradient descent, under the condition that Φ does not deviate too much from its linearization. In Sections 11.5 and 11.6, we then detail the implications for wide neural networks for two (slightly) different architectures. In particular, we will prove that gradient descent can find global minimizers when applied to (11.0.1b) for networks of very large width. We emphasize that this analysis treats the case of strong overparametrization, specifically when network width increases while keeping the number of data points m fixed.

11.1 Linear least-squares

Arguably one of the simplest machine learning algorithms is linear least-squares regression. Given data (11.0.1a), linear regression tries to fit a linear function $\Phi(\mathbf{x}, \mathbf{w}) := \mathbf{x}^\top \mathbf{w}$ in terms of \mathbf{w} by

minimizing $f(\mathbf{w})$ in (11.0.1b). With

$$\mathbf{A} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_m^\top \end{pmatrix} \in \mathbb{R}^{m \times d} \quad \text{and} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m \quad (11.1.1)$$

it holds

$$f(\mathbf{w}) = \|\mathbf{A}\mathbf{w} - \mathbf{y}\|^2. \quad (11.1.2)$$

Remark 11.1. More generally, the ansatz $\Phi(\mathbf{x}, (\mathbf{w}, b)) := \mathbf{w}^\top \mathbf{x} + b$ corresponds to

$$\Phi(\mathbf{x}, (\mathbf{w}, b)) = (1, \mathbf{x}^\top) \begin{pmatrix} b \\ \mathbf{w} \end{pmatrix}.$$

Therefore, additionally allowing for a bias can be treated analogously.

The model $\Phi(\mathbf{x}, \mathbf{w}) = \mathbf{x}^\top \mathbf{w}$ is linear in both \mathbf{x} and \mathbf{w} . In particular, $\mathbf{w} \mapsto f(\mathbf{w})$ is a convex function by Exercise 10.34, and we may apply the convergence results of Chapter 10 when using gradient based algorithms. If \mathbf{A} is invertible, then f has a unique minimizer given by $\mathbf{w}_* = \mathbf{A}^{-1}\mathbf{y}$. If $\text{rank}(\mathbf{A}) = d$, then f is strongly convex by Exercise 10.34, and there still exists a unique minimizer. If however $\text{rank}(\mathbf{A}) < d$, then $\ker(\mathbf{A}) \neq \{\mathbf{0}\}$ and there exist infinitely many minimizers of f . To ensure uniqueness, we look for the **minimum norm solution** (or minimum 2-norm solution)

$$\mathbf{w}_* := \operatorname{argmin}_{\{\mathbf{w} \in \mathbb{R}^d \mid f(\mathbf{w}) \leq f(\mathbf{v}) \ \forall \mathbf{v} \in \mathbb{R}^d\}} \|\mathbf{w}\|. \quad (11.1.3)$$

The following proposition establishes the uniqueness of \mathbf{w}_* and demonstrates that it can be represented as a superposition of the $(\mathbf{x}_i)_{i=1}^m$.

Proposition 11.2. *Let $\mathbf{A} \in \mathbb{R}^{m \times d}$ and $\mathbf{y} \in \mathbb{R}^m$ be as in (11.1.1). There exists a unique minimum 2-norm solution of (11.1.2). Denoting $\tilde{H} := \operatorname{span}\{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subseteq \mathbb{R}^d$, it is the unique element*

$$\mathbf{w}_* = \operatorname{argmin}_{\tilde{\mathbf{w}} \in \tilde{H}} f(\tilde{\mathbf{w}}) \in \tilde{H}. \quad (11.1.4)$$

Proof. We start with existence and uniqueness. Let $C \subseteq \mathbb{R}^m$ be the space spanned by the columns of \mathbf{A} . Then C is closed and convex, and therefore $\mathbf{y}_* = \operatorname{argmin}_{\tilde{\mathbf{y}} \in C} \|\mathbf{y} - \tilde{\mathbf{y}}\|$ exists and is unique (this is a fundamental property of Hilbert spaces, see Theorem B.14). In particular, the set $M = \{\mathbf{w} \in \mathbb{R}^d \mid \mathbf{A}\mathbf{w} = \mathbf{y}_*\} \subseteq \mathbb{R}^d$ of minimizers of f is not empty. Clearly M is also closed and convex. By the same argument as before, $\mathbf{w}_* = \operatorname{argmin}_{\mathbf{w}_* \in M} \|\mathbf{w}_*\|$ exists and is unique.

It remains to show (11.1.4). Denote by \mathbf{w}_* the minimum norm solution and decompose $\mathbf{w}_* = \tilde{\mathbf{w}} + \hat{\mathbf{w}}$ with $\tilde{\mathbf{w}} \in \tilde{H}$ and $\hat{\mathbf{w}} \in \tilde{H}^\perp$. We have $\mathbf{A}\mathbf{w}_* = \mathbf{A}\tilde{\mathbf{w}}$ and $\|\mathbf{w}_*\|^2 = \|\tilde{\mathbf{w}}\|^2 + \|\hat{\mathbf{w}}\|^2$. Since \mathbf{w}_* is the minimal norm solution it must hold $\hat{\mathbf{w}} = \mathbf{0}$. Thus $\mathbf{w}_* \in \tilde{H}$. Finally assume there exists a minimizer \mathbf{v} of f in \tilde{H} different from \mathbf{w}_* . Then $\mathbf{0} \neq \mathbf{w}_* - \mathbf{v} \in \tilde{H}$, and since \tilde{H} is spanned by the rows of \mathbf{A} we have $\mathbf{A}(\mathbf{w}_* - \mathbf{v}) \neq \mathbf{0}$. Thus $\mathbf{y}_* = \mathbf{A}\mathbf{w}_* \neq \mathbf{A}\mathbf{v}$, which contradicts that \mathbf{v} minimizes f . \square

The condition of minimizing the 2-norm is a form of regularization. Interestingly, gradient descent converges to the minimum norm solution for the quadratic objective (11.1.2), as long as \mathbf{w}_0 is initialized within $\tilde{H} = \text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ (e.g. $\mathbf{w}_0 = 0$). Therefore, it does not find an “arbitrary” minimizer but implicitly regularizes the problem in this sense. In the following $s_{\max}(\mathbf{A})$ denotes the maximal singular value of \mathbf{A} .

Theorem 11.3. *Let $\mathbf{A} \in \mathbb{R}^{m \times d}$ be as in (11.1.1), let $\mathbf{w}_0 = \tilde{\mathbf{w}}_0 + \hat{\mathbf{w}}_0$ where $\tilde{\mathbf{w}}_0 \in \tilde{H}$ and $\hat{\mathbf{w}}_0 \in \tilde{H}^\perp$. Fix $h \in (0, 1/(2s_{\max}(\mathbf{A})^2))$ and set*

$$\mathbf{w}_{k+1} := \mathbf{w}_k - h \nabla f(\mathbf{w}_k) \quad \text{for all } k \in \mathbb{N} \quad (11.1.5)$$

with f in (11.1.2). Then

$$\lim_{k \rightarrow \infty} \mathbf{w}_k = \mathbf{w}_* + \hat{\mathbf{w}}_0.$$

We sketch the argument in case $\mathbf{w}_0 \in \tilde{H}$, and leave the full proof to the reader, see Exercise 11.32. Note that \tilde{H} is the space spanned by the rows of \mathbf{A} (or the columns of \mathbf{A}^\top). The gradient of the objective function equals

$$\nabla f(\mathbf{w}) = 2\mathbf{A}^\top(\mathbf{A}\mathbf{w} - \mathbf{y}).$$

Therefore, if $\mathbf{w}_0 \in \tilde{H}$, then the iterates of gradient descent never leave the subspace \tilde{H} . By Exercise 10.34 and Theorem 10.11, for small enough step size, it holds $f(\mathbf{w}_k) \rightarrow 0$. By Proposition 11.2 there only exists one minimizer in \tilde{H} , corresponding to the minimum norm solution. Thus \mathbf{w}_k converges to the minimal norm solution.

11.2 Kernel least-squares

Let again $(\mathbf{x}_j, y_j) \in \mathbb{R}^d \times \mathbb{R}$, $j = 1, \dots, m$. In many applications linear models are too simplistic, and are not able to capture the true relation between \mathbf{x} and y . Kernel methods allow to overcome this problem by introducing nonlinearity in \mathbf{x} , but retaining linearity in the parameter \mathbf{w} .

Let H be a Hilbert space with inner product $\langle \cdot, \cdot \rangle_H$, that is also referred to as the **feature space**. For a (typically nonlinear) **feature map** $\phi : \mathbb{R}^d \rightarrow H$, consider the model

$$\Phi(\mathbf{x}, \mathbf{w}) = \langle \phi(\mathbf{x}), \mathbf{w} \rangle_H \quad (11.2.1)$$

with $\mathbf{w} \in H$. If $H = \mathbb{R}^n$, the components of ϕ are referred to as features. With the objective function

$$f(\mathbf{w}) := \sum_{j=1}^m (\langle \phi(\mathbf{x}_j), \mathbf{w} \rangle_H - y_j)^2 \quad \mathbf{w} \in H, \quad (11.2.2)$$

we wish to determine a minimizer of f . To ensure uniqueness and regularize the problem, we again consider the minimum H -norm solution

$$\mathbf{w}_* := \operatorname{argmin}_{\{\mathbf{w} \in H \mid f(\mathbf{w}) \leq f(\mathbf{v}) \ \forall \mathbf{v} \in H\}} \|\mathbf{w}\|_H.$$

As we will see below, \mathbf{w}_* is well-defined. We will call $\Phi(\mathbf{x}, \mathbf{w}_*) = \langle \phi(\mathbf{x}), \mathbf{w}_* \rangle_H$ the **kernel least squares estimator**. The nonlinearity of the feature map allows for more expressive models $\mathbf{x} \mapsto \Phi(\mathbf{x}, \mathbf{w})$ capable of capturing more complicated structures beyond linearity in the data.

Remark 11.4 (Gradient descent). Let $H = \mathbb{R}^n$ be equipped with the Euclidean inner product. Consider the sequence $(\mathbf{w}_k)_{k \in \mathbb{N}_0} \subseteq \mathbb{R}^n$ generated by gradient descent to minimize (11.2.2). Assuming sufficiently small step size, by Theorem 11.3 for $\mathbf{x} \in \mathbb{R}^d$

$$\lim_{k \rightarrow \infty} \Phi(\mathbf{x}, \mathbf{w}_k) = \langle \phi(\mathbf{x}), \mathbf{w}_* \rangle + \langle \phi(\mathbf{x}), \hat{\mathbf{w}}_0 \rangle. \quad (11.2.3)$$

Here, $\hat{\mathbf{w}}_0 \in \mathbb{R}^n$ denotes the orthogonal projection of $\mathbf{w}_0 \in \mathbb{R}^n$ onto \tilde{H}^\perp where $\tilde{H} := \text{span}\{\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_m)\}$. Gradient descent thus yields the kernel least squares estimator plus $\langle \phi(\mathbf{x}), \hat{\mathbf{w}}_0 \rangle$. Notably, on the set

$$\{\mathbf{x} \in \mathbb{R}^d \mid \phi(\mathbf{x}) \in \text{span}\{\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_m)\}\}, \quad (11.2.4)$$

(11.2.3) thus coincides with the kernel least squares estimator independent of the initialization \mathbf{w}_0 .

11.2.1 Examples

To motivate the concept of feature maps consider the following example from [153].

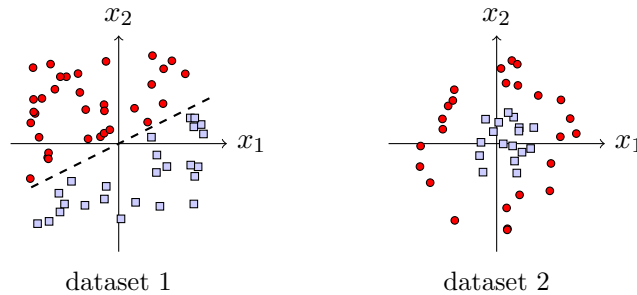
Example 11.5. Let $\mathbf{x}_i \in \mathbb{R}^2$ with associated labels $y_i \in \{-1, 1\}$ for $i = 1, \dots, m$. The goal is to find some model $\Phi(\cdot, \mathbf{w}) : \mathbb{R}^2 \rightarrow \mathbb{R}$, for which

$$\text{sign}(\Phi(\mathbf{x}, \mathbf{w})) \quad (11.2.5)$$

predicts the label y of \mathbf{x} . For a linear (in \mathbf{x}) model

$$\Phi(\mathbf{x}, (\mathbf{w}, b)) = \mathbf{x}^\top \mathbf{w} + b,$$

the decision boundary of (11.2.5) equals $\{\mathbf{x} \in \mathbb{R}^2 \mid \mathbf{x}^\top \mathbf{w} + b = 0\}$ in \mathbb{R}^2 . Hence, by adjusting \mathbf{w} and b , (11.2.5) can separate data by affine hyperplanes in \mathbb{R}^2 . Consider two datasets represented by light blue squares for $+1$ and red circles for -1 labels:

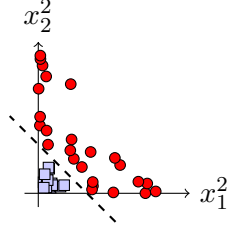


The first dataset is separable by an affine hyperplane as depicted by the dashed line. Thus a linear model is capable of correctly classifying all datapoints. For the second dataset this is not possible.

To enhance model expressivity, introduce a feature map $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^6$ via

$$\phi(\mathbf{x}) = (1, x_1, x_2, x_1 x_2, x_1^2, x_2^2)^\top \in \mathbb{R}^6 \quad \text{for all } \mathbf{x} \in \mathbb{R}^2. \quad (11.2.6)$$

For $\mathbf{w} \in \mathbb{R}^6$, this allows $\Phi(\mathbf{x}) = \mathbf{w}^\top \phi(\mathbf{x})$ to represent arbitrary polynomials of degree 2. With this kernel approach, the decision boundary of (11.2.5) becomes the set of all hyperplanes *in the feature space* passing through $\mathbf{0} \in \mathbb{R}^6$. Visualizing the last two features of the second dataset, we obtain



features 5 and 6 of dataset 2

Note how in the feature space \mathbb{R}^6 , the datapoints are again separated by such a hyperplane. Thus, with the feature map in (11.2.6), the predictor (11.2.5) can perfectly classify all points also for the second dataset.

In the above example we chose the feature space $H = \mathbb{R}^6$. It is also possible to work with infinite dimensional feature spaces as the next example demonstrates.

Example 11.6. Let $H = \ell^2(\mathbb{N})$ be the space of square summable sequences and $\phi : \mathbb{R}^d \rightarrow \ell^2(\mathbb{N})$ some map. Fitting the corresponding model

$$\Phi(\mathbf{x}, \mathbf{w}) = \langle \phi(\mathbf{x}), \mathbf{w} \rangle_{\ell^2} = \sum_{i \in \mathbb{N}} \phi_i(\mathbf{x}) w_i$$

to data $(\mathbf{x}_i, y_i)_{i=1}^m$ requires to minimize

$$f(\mathbf{w}) = \sum_{j=1}^m \left(\left(\sum_{i \in \mathbb{N}} \phi_i(\mathbf{x}_j) w_i \right) - y_j \right)^2 \quad \mathbf{w} \in \ell^2(\mathbb{N}).$$

Hence we have to determine an *infinite sequence* of parameters $(w_i)_{i \in \mathbb{N}}$.

11.2.2 Kernel trick

At first glance, computing a (minimal H -norm) minimizer \mathbf{w} in the possibly infinite-dimensional Hilbert space H seems infeasible. The so-called *kernel trick* allows to do this computation. To explain it, we first revisit the foundational representer theorem.

Theorem 11.7 (Representer theorem). *There is a unique minimum H -norm solution $\mathbf{w}_* \in H$ of (11.2.2). With $\tilde{H} := \text{span}\{\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_m)\}$ it equals the unique element*

$$\mathbf{w}_* = \operatorname{argmin}_{\tilde{\mathbf{w}} \in \tilde{H}} f(\tilde{\mathbf{w}}) \in \tilde{H}. \quad (11.2.7)$$

Proof. Let $\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_n$ be a basis of \tilde{H} . If $\tilde{H} = \{0\}$ the statement is trivial, so we assume $1 \leq n \leq m$. Let $\mathbf{A} = (\langle \phi(\mathbf{x}_i), \tilde{\mathbf{w}}_j \rangle)_{ij} \in \mathbb{R}^{m \times n}$. Every $\tilde{\mathbf{w}} \in \tilde{H}$ has a unique representation $\tilde{\mathbf{w}} = \sum_{j=1}^n \alpha_j \tilde{\mathbf{w}}_j$ for some $\boldsymbol{\alpha} \in \mathbb{R}^n$. With this ansatz

$$f(\tilde{\mathbf{w}}) = \sum_{i=1}^m (\langle \phi(\mathbf{x}_i), \tilde{\mathbf{w}} \rangle - y_i)^2 = \sum_{i=1}^m \left(\sum_{j=1}^n \langle \phi(\mathbf{x}_i), \tilde{\mathbf{w}}_j \rangle \alpha_j - y_i \right)^2 = \|\mathbf{A}\boldsymbol{\alpha} - \mathbf{y}\|^2. \quad (11.2.8)$$

Note that $\mathbf{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is injective since for every $\boldsymbol{\alpha} \in \mathbb{R}^n \setminus \{0\}$ holds $\sum_{j=1}^n \alpha_j \tilde{\mathbf{w}}_j \in \tilde{H} \setminus \{0\}$ and hence $\mathbf{A}\boldsymbol{\alpha} = (\langle \phi(\mathbf{x}_i), \sum_{j=1}^n \alpha_j \tilde{\mathbf{w}}_j \rangle)_{i=1}^m \neq 0$. Therefore, there exists a unique minimizer $\boldsymbol{\alpha} \in \mathbb{R}^n$ of the right-hand side of (11.2.8), and thus there exists a unique minimizer $\mathbf{w}_* \in \tilde{H}$ in (11.2.7).

For arbitrary $\mathbf{w} \in H$ we wish to show $f(\mathbf{w}) \geq f(\mathbf{w}_*)$, so that \mathbf{w}_* minimizes f in H . Decompose $\mathbf{w} = \tilde{\mathbf{w}} + \hat{\mathbf{w}}$ with $\tilde{\mathbf{w}} \in \tilde{H}$ and $\hat{\mathbf{w}} \in \tilde{H}^\perp$, i.e. $\langle \phi(\mathbf{x}_j), \hat{\mathbf{w}} \rangle_H = 0$ for all $j = 1, \dots, m$. Then, using that \mathbf{w}_* minimizes f in \tilde{H} ,

$$f(\mathbf{w}) = \sum_{j=1}^m (\langle \phi(\mathbf{x}_j), \mathbf{w} \rangle_H - y_j)^2 = \sum_{j=1}^m (\langle \phi(\mathbf{x}_j), \tilde{\mathbf{w}} \rangle_H - y_j)^2 = f(\tilde{\mathbf{w}}) \geq f(\mathbf{w}_*).$$

Finally, let $\mathbf{w} \in H$ be any minimizer of f in H different from \mathbf{w}_* . It remains to show $\|\mathbf{w}\|_H > \|\mathbf{w}_*\|_H$. Decompose again $\mathbf{w} = \tilde{\mathbf{w}} + \hat{\mathbf{w}}$ with $\tilde{\mathbf{w}} \in \tilde{H}$ and $\hat{\mathbf{w}} \in \tilde{H}^\perp$. As above $f(\mathbf{w}) = f(\tilde{\mathbf{w}})$ and thus $\tilde{\mathbf{w}}$ is a minimizer of f . Uniqueness of \mathbf{w}_* in (11.2.7) implies $\tilde{\mathbf{w}} = \mathbf{w}_*$. Therefore $\hat{\mathbf{w}} \neq 0$ and $\|\mathbf{w}_*\|_H^2 < \|\tilde{\mathbf{w}}\|_H^2 + \|\hat{\mathbf{w}}\|_H^2 = \|\mathbf{w}\|_H^2$. \square

Instead of looking for the minimum norm minimizer \mathbf{w}_* in the Hilbert space H , by Proposition 11.2 it suffices to determine the unique minimizer in the at most m -dimensional subspace \tilde{H} spanned by $\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_m)$. This significantly simplifies the problem. To do so we first introduce the notion of kernels.

Definition 11.8. A symmetric function $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is called a **kernel** if for any $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^d$ the **kernel matrix** $\mathbf{G} = (K(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^m \in \mathbb{R}^{m \times m}$ is symmetric positive semidefinite.

Given a feature map $\phi : \mathbb{R}^d \rightarrow H$, it is easy to check that

$$K(\mathbf{x}, \mathbf{x}') := \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_H \quad \text{for all } \mathbf{x}, \mathbf{x}' \in \mathbb{R}^d,$$

defines a kernel. The corresponding kernel matrix $\mathbf{G} \in \mathbb{R}^{m \times m}$ is given by

$$G_{ij} = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle_H = K(\mathbf{x}_i, \mathbf{x}_j).$$

With the ansatz $\mathbf{w} = \sum_{j=1}^m \alpha_j \phi(\mathbf{x}_j)$, minimizing the objective (11.2.2) in \tilde{H} is equivalent to minimizing

$$\|\mathbf{G}\boldsymbol{\alpha} - \mathbf{y}\|^2, \tag{11.2.9}$$

in $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_m) \in \mathbb{R}^m$.

Proposition 11.9. Let $\boldsymbol{\alpha} \in \mathbb{R}^m$ be any minimizer of (11.2.9). Then $\mathbf{w}_* = \sum_{j=1}^m \alpha_j \phi(\mathbf{x}_j)$ is the unique minimum H -norm solution of (11.2.2).

Proposition 11.9, the proof of which is left as an exercise, suggests the following algorithm to compute the kernel least squares estimator:

- (i) compute the kernel matrix $\mathbf{G} = (K(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^m$,
- (ii) determine a minimizer $\boldsymbol{\alpha} \in \mathbb{R}^m$ of $\|\mathbf{G}\boldsymbol{\alpha} - \mathbf{y}\|$,
- (iii) evaluate $\Phi(\mathbf{x}, \mathbf{w}_*)$ via

$$\Phi(\mathbf{x}, \mathbf{w}_*) = \left\langle \phi(\mathbf{x}), \sum_{j=1}^m \alpha_j \phi(\mathbf{x}_j) \right\rangle_H = \sum_{j=1}^m \alpha_j K(\mathbf{x}, \mathbf{x}_j). \quad (11.2.10)$$

Thus, minimizing (11.2.2) and expressing the kernel least squares estimator does neither require explicit knowledge of the feature map ϕ nor of the minimum norm solution $\mathbf{w}_* \in H$. It is sufficient to choose a kernel map $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$; this is known as the kernel trick. Given a kernel K , we will therefore also refer to (11.2.10) as the kernel least squares estimator without specifying H or ϕ .

Example 11.10. Common examples of kernels include the **polynomial kernel**

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + c)^r \quad c \geq 0, \quad r \in \mathbb{N},$$

the radial basis function (**RBF**) **kernel**

$$K(\mathbf{x}, \mathbf{x}') = \exp(-c\|\mathbf{x} - \mathbf{x}'\|^2) \quad c > 0,$$

and the **Laplace kernel**

$$K(\mathbf{x}, \mathbf{x}') = \exp(-c\|\mathbf{x} - \mathbf{x}'\|) \quad c > 0.$$

Remark 11.11. If $\Omega \subseteq \mathbb{R}^d$ is compact and $K : \Omega \times \Omega \rightarrow \mathbb{R}$ is a continuous kernel, then Mercer's theorem implies existence of a Hilbert space H and a feature map $\phi : \mathbb{R}^d \rightarrow H$ such that

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_H \quad \text{for all } \mathbf{x}, \mathbf{x}' \in \Omega,$$

i.e. K is the corresponding kernel. See for instance [216, Thm. 4.49].

11.3 Tangent kernel

Consider again a general model $\Phi(\mathbf{x}, \mathbf{w})$ with input $\mathbf{x} \in \mathbb{R}^d$ and parameters $\mathbf{w} \in \mathbb{R}^n$. The goal remains to minimize the square loss objective (11.0.1b) given the data (11.0.1a). If $\mathbf{w} \mapsto \Phi(\mathbf{x}, \mathbf{w})$ is not linear, then unlike in Sections 11.1 and 11.2, the objective function (11.0.1b) is in general not convex, and most results on first order methods in Chapter 10 are not directly applicable.

We now simplify the situation by *linearizing the model in $\mathbf{w} \in \mathbb{R}^n$ around the initialization*: Fixing $\mathbf{w}_0 \in \mathbb{R}^n$, let

$$\Phi^{\text{lin}}(\mathbf{x}, \mathbf{w}) := \Phi(\mathbf{x}, \mathbf{w}_0) + \nabla_{\mathbf{w}} \Phi(\mathbf{x}, \mathbf{w}_0)^\top (\mathbf{w} - \mathbf{w}_0) \quad \text{for all } \mathbf{w} \in \mathbb{R}^n, \quad (11.3.1)$$

which is the first order Taylor approximation of Φ around the initial parameter \mathbf{w}_0 . Introduce the notation

$$\delta_i := \Phi(\mathbf{x}_i, \mathbf{w}_0) - \nabla_{\mathbf{w}} \Phi(\mathbf{x}_i, \mathbf{w}_0)^\top \mathbf{w}_0 - y_i \quad \text{for all } i = 1, \dots, m. \quad (11.3.2)$$

The square loss for the linearized model then reads

$$\begin{aligned} f^{\text{lin}}(\mathbf{w}) &:= \sum_{j=1}^m (\Phi^{\text{lin}}(\mathbf{x}_j, \mathbf{w}) - y_j)^2 \\ &= \sum_{j=1}^m (\langle \nabla_{\mathbf{w}} \Phi(\mathbf{x}_j, \mathbf{w}_0), \mathbf{w} \rangle + \delta_j)^2, \end{aligned} \quad (11.3.3)$$

where $\langle \cdot, \cdot \rangle$ stands for the Euclidean inner product in \mathbb{R}^n . Comparing with (11.2.2), minimizing f^{lin} corresponds to a kernel least squares regression with feature map

$$\phi(\mathbf{x}) = \nabla_{\mathbf{w}} \Phi(\mathbf{x}, \mathbf{w}_0) \in \mathbb{R}^n.$$

The corresponding kernel is

$$\hat{K}_n(\mathbf{x}, \mathbf{x}') = \langle \nabla_{\mathbf{w}} \Phi(\mathbf{x}, \mathbf{w}_0), \nabla_{\mathbf{w}} \Phi(\mathbf{x}', \mathbf{w}_0) \rangle. \quad (11.3.4)$$

We refer to \hat{K}_n as the empirical **tangent kernel**, as it arises from the first order Taylor approximation (the tangent) of the original model Φ around initialization \mathbf{w}_0 . Note that the kernel depends on the choice of \mathbf{w}_0 . As explained in Remark 11.4, training Φ^{lin} with gradient descent yields the kernel least-squares estimator with kernel \hat{K}_n plus an additional term depending on \mathbf{w}_0 .

Of course the linearized model Φ^{lin} only captures the behaviour of Φ for parameters \mathbf{w} that are close to \mathbf{w}_0 . If we assume for the moment that during training of Φ , the parameters remain close to initialization, then we can expect similar behaviour and performance of Φ and Φ^{lin} . Under certain assumptions, we will see in the next sections that this is precisely what happens, when the width of a neural network increases. Before we make this precise, in Section 11.4 we investigate whether gradient descent applied to $f(\mathbf{w})$ will find a *global* minimizer, under the assumption that Φ^{lin} is a good approximation of Φ .

11.4 Convergence to global minimizers

Intuitively, if $\mathbf{w} \mapsto \Phi(\mathbf{x}, \mathbf{w})$ is not linear but “close enough to its linearization” Φ^{lin} defined in (11.3.1), we expect that the objective function is close to a convex function and gradient descent can still find global minimizers of (11.0.1b). To motivate this, consider Figures 11.1 and 11.2 where we chose the number of training data $m = 1$ and the number of parameters $n = 1$. As we can see, essentially we require the difference of Φ and Φ^{lin} and of their derivatives to be small in a neighbourhood of \mathbf{w}_0 . The size of the neighbourhood crucially depends on the initial error $\Phi(\mathbf{x}_1, \mathbf{w}_0) - y_1$, and on the size of the derivative $\frac{d}{d\mathbf{w}} \Phi(\mathbf{x}_1, \mathbf{w}_0)$.

For general m and n , we now make the required assumptions on Φ precise.

Assumption 11.12. Let $\Phi \in C^1(\mathbb{R}^d \times \mathbb{R}^n)$ and $\mathbf{w}_0 \in \mathbb{R}^n$. There exist constants $r > 0$, $U, L < \infty$ and $0 < \lambda_{\min} \leq \lambda_{\max} < \infty$ such that

- (a) the kernel matrix of the empirical tangent kernel

$$(\hat{K}_n(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^m = (\langle \nabla_{\mathbf{w}} \Phi(\mathbf{x}_i, \mathbf{w}_0), \nabla_{\mathbf{w}} \Phi(\mathbf{x}_j, \mathbf{w}_0) \rangle)_{i,j=1}^m \in \mathbb{R}^{m \times m} \quad (11.4.1)$$

is regular and its eigenvalues belong to $[\lambda_{\min}, \lambda_{\max}]$,

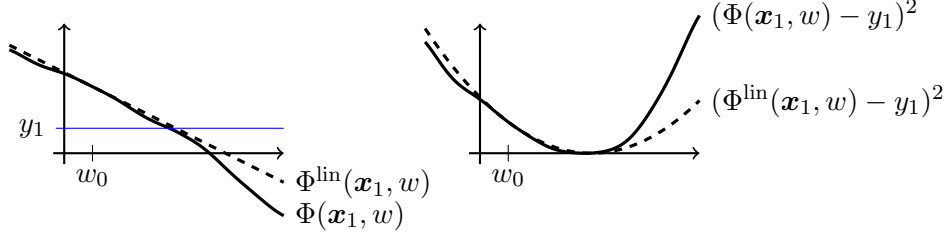


Figure 11.1: Graph of a model $w \mapsto \Phi(\mathbf{x}_1, w)$ and its linearization $w \mapsto \Phi^{\text{lin}}(\mathbf{x}_1, w)$ at the initial parameter w_0 , s.t. $\frac{d}{dw}\Phi(\mathbf{x}_1, w_0) \neq 0$. If Φ and Φ^{lin} are close, then there exists w s.t. $\Phi(\mathbf{x}_1, w) = y_1$ (left). If the derivatives are also close, the loss $(\Phi(\mathbf{x}_1, w) - y_1)^2$ is nearly convex in w , and gradient descent finds a global minimizer (right).

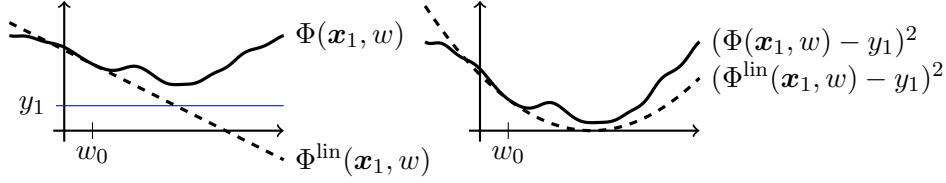


Figure 11.2: Same as Figure 11.1. If Φ and Φ^{lin} are not close, there need not exist w such that $\Phi(\mathbf{x}_1, w) = y_1$, and gradient descent need not converge to a global minimizer.

(b) for all $i \in \{1, \dots, m\}$ holds

$$\begin{aligned} \|\nabla_{\mathbf{w}}\Phi(\mathbf{x}_i, \mathbf{w})\| &\leq U && \text{for all } \mathbf{w} \in B_r(\mathbf{w}_0) \\ \|\nabla_{\mathbf{w}}\Phi(\mathbf{x}_i, \mathbf{w}) - \nabla_{\mathbf{w}}\Phi(\mathbf{x}_i, \mathbf{v})\| &\leq L\|\mathbf{w} - \mathbf{v}\| && \text{for all } \mathbf{w}, \mathbf{v} \in B_r(\mathbf{w}_0), \end{aligned} \quad (11.4.2)$$

(c) and

$$L \leq \frac{\lambda_{\min}^2}{12m^{3/2}U^2\sqrt{f(\mathbf{w}_0)}} \quad \text{and} \quad r = \frac{2\sqrt{m}U}{\lambda_{\min}}\sqrt{f(\mathbf{w}_0)}. \quad (11.4.3)$$

The regularity of the kernel matrix in Assumption 11.12 (a) is equivalent to $(\nabla_{\mathbf{w}}\Phi(\mathbf{x}_i, \mathbf{w}_0))_{i=1}^m \in \mathbb{R}^{m \times n}$ having full rank $m \leq n$ (in particular we have at least as many parameters n as training data m). In the context of Figure 11.1, this means that $\frac{d}{dw}\Phi(\mathbf{x}_1, w_0) \neq 0$ and thus Φ^{lin} is not a constant function. This condition guarantees that there exists \mathbf{w} such that $\Phi^{\text{lin}}(\mathbf{x}_i, \mathbf{w}) = y_i$ for all $i = 1, \dots, m$. In other words, already the linearized model Φ^{lin} is sufficiently expressive to interpolate the data. Assumption 11.12 (b) formalizes the closeness condition of Φ and Φ^{lin} . Apart from giving an upper bound on $\nabla_{\mathbf{w}}\Phi(\mathbf{x}_i, \mathbf{w})$, it assumes $\mathbf{w} \mapsto \Phi(\mathbf{x}_i, \mathbf{w})$ to be L -smooth in a ball of radius $r > 0$ around \mathbf{w}_0 , for all $i = 1, \dots, m$. This allows to control how far $\Phi(\mathbf{x}_i, \mathbf{w})$ and $\Phi^{\text{lin}}(\mathbf{x}_i, \mathbf{w})$ and their derivatives may deviate from each other for \mathbf{w} in this ball. Finally Assumption 11.12 (c) ties together all constants, ensuring the full model to be sufficiently close to its linearization in a large enough neighbourhood of \mathbf{w}_0 .

We are now ready to state the following theorem, which is a variant of [129, Thm. G.1]. In Section 11.5 we will see that its main requirement—Assumption 11.12—is satisfied with high probability for certain (wide) neural networks.

Theorem 11.13. *Let Assumption 11.12 be satisfied and fix a positive learning rate*

$$h \leq \frac{1}{\lambda_{\min} + \lambda_{\max}}. \quad (11.4.4)$$

Set for all $k \in \mathbb{N}$

$$\mathbf{w}_{k+1} = \mathbf{w}_k - h \nabla f(\mathbf{w}_k). \quad (11.4.5)$$

It then holds for all $k \in \mathbb{N}$

$$\|\mathbf{w}_k - \mathbf{w}_0\| \leq \frac{2\sqrt{m}U}{\lambda_{\min}} \sqrt{f(\mathbf{w}_0)} \quad (11.4.6a)$$

$$f(\mathbf{w}_k) \leq (1 - h\lambda_{\min})^{2k} f(\mathbf{w}_0). \quad (11.4.6b)$$

Proof. In the following denote the error in prediction by

$$E(\mathbf{w}) := (\Phi(\mathbf{x}_i, \mathbf{w}) - y_i)_{i=1}^m \in \mathbb{R}^m$$

such that

$$\nabla E(\mathbf{w}) = (\nabla_{\mathbf{w}} \Phi(\mathbf{x}_i, \mathbf{w}))_{i=1}^m \in \mathbb{R}^{m \times n}$$

and with the empirical tangent kernel \hat{K}_n in Assumption 11.12

$$\nabla E(\mathbf{w}) \nabla E(\mathbf{w})^\top = (\hat{K}_n(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^m \in \mathbb{R}^{m \times m}. \quad (11.4.7)$$

Moreover, (11.4.2) gives

$$\|\nabla E(\mathbf{w})\|^2 \leq \|\nabla E(\mathbf{w})\|_F^2 = \sum_{i=1}^m \|\nabla \Phi(\mathbf{x}_i, \mathbf{w})\|^2 \leq mU^2 \quad \text{for all } \mathbf{w} \in B_r(\mathbf{w}_0), \quad (11.4.8a)$$

and similarly

$$\begin{aligned} \|\nabla E(\mathbf{w}) - \nabla E(\mathbf{v})\|^2 &\leq \sum_{i=1}^m \|\nabla_{\mathbf{w}} \Phi(\mathbf{x}_i, \mathbf{w}) - \nabla_{\mathbf{w}} \Phi(\mathbf{x}_i, \mathbf{v})\|^2 \\ &\leq mL^2 \|\mathbf{w} - \mathbf{v}\|^2 \quad \text{for all } \mathbf{w}, \mathbf{v} \in B_r(\mathbf{w}_0). \end{aligned} \quad (11.4.8b)$$

Denote $c := 1 - h\lambda_{\min} \in (0, 1)$. We use induction over k to prove

$$\sum_{j=0}^{k-1} \|\mathbf{w}_{j+1} - \mathbf{w}_j\| \leq h2\sqrt{m}U \|E(\mathbf{w}_0)\| \sum_{j=0}^{k-1} c^j, \quad (11.4.9a)$$

$$\|E(\mathbf{w}_k)\|^2 \leq \|E(\mathbf{w}_0)\|^2 c^{2k}, \quad (11.4.9b)$$

for all $k \in \mathbb{N}_0$ and where an empty sum is understood as zero. Since $\sum_{j=0}^{\infty} c^j = (1-c)^{-1} = (h\lambda_{\min})^{-1}$ and $f(\mathbf{w}_k) = \|E(\mathbf{w}_k)\|^2$, these inequalities directly imply (11.4.6).

The case $k = 0$ is trivial. For the induction step, assume (11.4.9) holds for some $k \in \mathbb{N}_0$.

Step 1. We show (11.4.9a) for $k + 1$. The induction assumption and (11.4.3) give

$$\|\mathbf{w}_k - \mathbf{w}_0\| \leq 2h\sqrt{m}U\|E(\mathbf{w}_0)\| \sum_{j=0}^{\infty} c^j = \frac{2\sqrt{m}U}{\lambda_{\min}} \sqrt{f(\mathbf{w}_0)} = r, \quad (11.4.10)$$

and thus $\mathbf{w}_k \in B_r(\mathbf{w}_0)$. Next

$$\nabla f(\mathbf{w}_k) = \nabla(E(\mathbf{w}_k)^\top E(\mathbf{w}_k)) = 2\nabla E(\mathbf{w}_k)^\top E(\mathbf{w}_k). \quad (11.4.11)$$

Using the iteration rule (11.4.5), the bound (11.4.8a), and (11.4.9b)

$$\begin{aligned} \|\mathbf{w}_{k+1} - \mathbf{w}_k\| &= 2h\|\nabla E(\mathbf{w}_k)^\top E(\mathbf{w}_k)\| \\ &\leq 2h\sqrt{m}U\|E(\mathbf{w}_k)\| \\ &\leq 2h\sqrt{m}U\|E(\mathbf{w}_0)\|c^k. \end{aligned}$$

This shows (11.4.9a) for $k + 1$. In particular, as in (11.4.10) we conclude

$$\mathbf{w}_{k+1}, \mathbf{w}_k \in B_r(\mathbf{w}_0). \quad (11.4.12)$$

Step 2. We show (11.4.9b) for $k + 1$. Since E is continuously differentiable, there exists $\tilde{\mathbf{w}}_k$ in the convex hull of \mathbf{w}_k and \mathbf{w}_{k+1} such that

$$E(\mathbf{w}_{k+1}) = E(\mathbf{w}_k) + \nabla E(\tilde{\mathbf{w}}_k)(\mathbf{w}_{k+1} - \mathbf{w}_k) = E(\mathbf{w}_k) - h\nabla E(\tilde{\mathbf{w}}_k)\nabla f(\mathbf{w}_k),$$

and thus by (11.4.11)

$$\begin{aligned} E(\mathbf{w}_{k+1}) &= E(\mathbf{w}_k) - 2h\nabla E(\tilde{\mathbf{w}}_k)\nabla E(\mathbf{w}_k)^\top E(\mathbf{w}_k) \\ &= (\mathbf{I}_m - 2h\nabla E(\tilde{\mathbf{w}}_k)\nabla E(\mathbf{w}_k)^\top)E(\mathbf{w}_k), \end{aligned}$$

where $\mathbf{I}_m \in \mathbb{R}^{m \times m}$ is the identity matrix. We wish to show that

$$\|\mathbf{I}_m - 2h\nabla E(\tilde{\mathbf{w}}_k)\nabla E(\mathbf{w}_k)^\top\| \leq c, \quad (11.4.13)$$

which then implies (11.4.9b) for $k + 1$ and concludes the proof.

Using (11.4.8) and the fact that $\mathbf{w}_k, \tilde{\mathbf{w}}_k \in B_r(\mathbf{w}_0)$ by (11.4.12),

$$\begin{aligned} &\|\nabla E(\tilde{\mathbf{w}}_k)\nabla E(\mathbf{w}_k)^\top - \nabla E(\mathbf{w}_0)\nabla E(\mathbf{w}_0)^\top\| \\ &\leq \|\nabla E(\tilde{\mathbf{w}}_k)\nabla E(\mathbf{w}_k)^\top - \nabla E(\mathbf{w}_k)\nabla E(\mathbf{w}_k)^\top\| \\ &\quad + \|\nabla E(\tilde{\mathbf{w}}_k)\nabla E(\mathbf{w}_k)^\top - \nabla E(\mathbf{w}_k)\nabla E(\mathbf{w}_0)^\top\| \\ &\quad + \|\nabla E(\tilde{\mathbf{w}}_k)\nabla E(\mathbf{w}_0)^\top - \nabla E(\mathbf{w}_0)\nabla E(\mathbf{w}_0)^\top\| \\ &\leq 3mULr. \end{aligned}$$

Since the eigenvalues of $\nabla E(\mathbf{w}_0)\nabla E(\mathbf{w}_0)^\top$ belong to $[\lambda_{\min}, \lambda_{\max}]$ by (11.4.7) and Assumption 11.12 (a), as long as $h \leq (\lambda_{\min} + \lambda_{\max})^{-1}$ we have

$$\begin{aligned} \|\mathbf{I}_m - 2h\nabla E(\tilde{\mathbf{w}}_k)\nabla E(\mathbf{w}_k)^\top\| &\leq \|\mathbf{I}_m - 2h\nabla E(\mathbf{w}_0)\nabla E(\mathbf{w}_0)^\top\| + 6hmULr \\ &\leq 1 - 2h\lambda_{\min} + 6hmULr \\ &\leq 1 - 2h(\lambda_{\min} - 3mULr) \\ &\leq 1 - h\lambda_{\min} = c, \end{aligned}$$

where we have used the equality for r and the upper bound for L in (11.4.3). \square

Let us emphasize the main statement of Theorem 11.13. By (11.4.6b), full batch gradient descent (11.4.5) achieves zero loss in the limit, i.e. the data is interpolated by the limiting model. In particular, this yields convergence for the (possibly nonconvex) optimization problem of minimizing $f(\mathbf{w})$.

11.5 Training dynamics for LeCun initialization

In this and the next section we discuss the implications of Theorem 11.13 for wide neural networks. For ease of presentation we focus on shallow networks with only one hidden layer, but stress that similar considerations also hold for deep networks, see the bibliography section.

11.5.1 Architecture

Let $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}$ be a neural network of depth one and width $n \in \mathbb{N}$ of type

$$\Phi(\mathbf{x}, \mathbf{w}) = \mathbf{v}^\top \sigma(\mathbf{U}\mathbf{x} + \mathbf{b}) + c. \quad (11.5.1)$$

Here $\mathbf{x} \in \mathbb{R}^d$ is the input, and $\mathbf{U} \in \mathbb{R}^{n \times d}$, $\mathbf{v} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^n$ and $c \in \mathbb{R}$ are the parameters which we collect in the vector $\mathbf{w} = (\mathbf{U}, \mathbf{b}, \mathbf{v}, c) \in \mathbb{R}^{n(d+2)+1}$ (with \mathbf{U} suitably reshaped). For future reference we note that

$$\begin{aligned} \nabla_{\mathbf{U}} \Phi(\mathbf{x}, \mathbf{w}) &= (\mathbf{v} \odot \sigma'(\mathbf{U}\mathbf{x} + \mathbf{b})) \mathbf{x}^\top \in \mathbb{R}^{n \times d} \\ \nabla_{\mathbf{b}} \Phi(\mathbf{x}, \mathbf{w}) &= \mathbf{v} \odot \sigma'(\mathbf{U}\mathbf{x} + \mathbf{b}) \in \mathbb{R}^n \\ \nabla_{\mathbf{v}} \Phi(\mathbf{x}, \mathbf{w}) &= \sigma(\mathbf{U}\mathbf{x} + \mathbf{b}) \in \mathbb{R}^n \\ \nabla_c \Phi(\mathbf{x}, \mathbf{w}) &= 1 \in \mathbb{R}, \end{aligned} \quad (11.5.2)$$

where \odot denotes the Hadamard product. We also write $\nabla_{\mathbf{w}} \Phi(\mathbf{x}, \mathbf{w}) \in \mathbb{R}^{n(d+2)+1}$ to denote the full gradient with respect to all parameters.

In practice, it is common to initialize the weights randomly, and in this section we consider so-called LeCun initialization. The following condition on the distribution used for this initialization will be assumed throughout the rest of Section 11.5.

Assumption 11.14. The distribution \mathcal{D} on \mathbb{R} has expectation zero, variance one, and finite moments up to order eight.

To explicitly indicate the expectation and variance in the notation, we also write $\mathcal{D}(0, 1)$ instead of \mathcal{D} , and for $\mu \in \mathbb{R}$ and $\varsigma > 0$ we use $\mathcal{D}(\mu, \varsigma^2)$ to denote the corresponding scaled and shifted measure with expectation μ and variance ς^2 ; thus, if $X \sim \mathcal{D}(0, 1)$ then $\mu + \varsigma X \sim \mathcal{D}(\mu, \varsigma^2)$. LeCun initialization [127] sets the variance of the weights in each layer to be reciprocal to the input dimension of the layer, thereby normalizing the output variance across all network nodes. The initial parameters

$$\mathbf{w}_0 = (\mathbf{U}_0, \mathbf{b}_0, \mathbf{v}_0, c_0)$$

are thus randomly initialized with components

$$U_{0;ij} \stackrel{\text{iid}}{\sim} \mathcal{D}\left(0, \frac{1}{d}\right), \quad v_{0;i} \stackrel{\text{iid}}{\sim} \mathcal{D}\left(0, \frac{1}{n}\right), \quad b_{0;i}, c_0 = 0, \quad (11.5.3)$$

independently for all $i = 1, \dots, n$, $j = 1, \dots, d$. For a fixed $\varsigma > 0$ one might choose variances ς^2/d and ς^2/n in (11.5.3), which would require only minor modifications in the rest of this section. Biases

are set to zero for simplicity, with nonzero initialization discussed in the exercises. All expectations and probabilities in Section 11.5 are understood with respect to this random initialization.

Example 11.15. Typical examples for $\mathcal{D}(0, 1)$ are the standard normal distribution on \mathbb{R} or the uniform distribution on $[-\sqrt{3}, \sqrt{3}]$.

11.5.2 Neural tangent kernel

We begin our analysis by investigating the empirical tangent kernel

$$\hat{K}_n(\mathbf{x}, \mathbf{z}) = \langle \nabla_{\mathbf{w}} \Phi(\mathbf{x}, \mathbf{w}_0), \nabla_{\mathbf{w}} \Phi(\mathbf{z}, \mathbf{w}_0) \rangle$$

of the shallow network (11.5.1). Scaled properly, it converges in the infinite width limit $n \rightarrow \infty$ towards a specific kernel known as the **neural tangent kernel** (NTK). Its precise formula depends on the architecture and initialization. For the LeCun initialization (11.5.3) we denote it by K^{LC} .

Theorem 11.16. *Let $R < \infty$ such that $|\sigma(x)| \leq R \cdot (1 + |x|)$ and $|\sigma'(x)| \leq R \cdot (1 + |x|)$ for all $x \in \mathbb{R}$. For any $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$ and $u_i \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1/d)$, $i = 1, \dots, d$, it then holds*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \hat{K}_n(\mathbf{x}, \mathbf{z}) = \mathbb{E}[\sigma(\mathbf{u}^\top \mathbf{x}) \sigma(\mathbf{u}^\top \mathbf{z})] =: K^{\text{LC}}(\mathbf{x}, \mathbf{z})$$

almost surely.

Moreover, for every $\delta, \varepsilon > 0$ there exists $n_0(\delta, \varepsilon, R) \in \mathbb{N}$ such that for all $n \geq n_0$ and all $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$ with $\|\mathbf{x}\|, \|\mathbf{z}\| \leq R$

$$\mathbb{P} \left[\left\| \frac{1}{n} \hat{K}_n(\mathbf{x}, \mathbf{z}) - K^{\text{LC}}(\mathbf{x}, \mathbf{z}) \right\| < \varepsilon \right] \geq 1 - \delta.$$

Proof. Denote $\mathbf{x}^{(1)} = \mathbf{U}_0 \mathbf{x} + \mathbf{b}_0 \in \mathbb{R}^n$ and $\mathbf{z}^{(1)} = \mathbf{U}_0 \mathbf{z} + \mathbf{b}_0 \in \mathbb{R}^n$. Due to the initialization (11.5.3) and our assumptions on $\mathcal{D}(0, 1)$, the components

$$x_i^{(1)} = \sum_{j=1}^d U_{0;ij} x_j \sim \mathbf{u}^\top \mathbf{x} \quad i = 1, \dots, n$$

are i.i.d. with finite p th moment (independent of n) for all $1 \leq p \leq 8$. Due to the linear growth bound on σ and σ' , the same holds for the $(\sigma(x_i^{(1)}))_{i=1}^n$ and the $(\sigma'(x_i^{(1)}))_{i=1}^n$. Similarly, the $(\sigma(z_i^{(1)}))_{i=1}^n$ and $(\sigma'(z_i^{(1)}))_{i=1}^n$ are collections of i.i.d. random variables with finite p th moment for all $1 \leq p \leq 8$.

Denote $\tilde{v}_i = \sqrt{n} v_{0;i}$ such that $\tilde{v}_i \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1)$. By (11.5.2)

$$\frac{1}{n} \hat{K}_n(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^\top \mathbf{z}) \frac{1}{n^2} \sum_{i=1}^n \tilde{v}_i^2 \sigma'(x_i^{(1)}) \sigma'(z_i^{(1)}) + \frac{1}{n} \sum_{i=1}^n \sigma(x_i^{(1)}) \sigma(z_i^{(1)}) + \frac{1}{n}.$$

Since

$$\frac{1}{n} \sum_{i=1}^n \tilde{v}_i^2 \sigma'(x_i^{(1)}) \sigma'(z_i^{(1)}) \tag{11.5.4}$$

is an average over i.i.d. random variables with finite variance, the law of large numbers implies almost sure convergence of this expression towards

$$\begin{aligned}\mathbb{E}[\tilde{v}_i^2 \sigma'(x_i^{(1)}) \sigma'(z_i^{(1)})] &= \mathbb{E}[\tilde{v}_i^2] \mathbb{E}[\sigma'(x_i^{(1)}) \sigma'(z_i^{(1)})] \\ &= \mathbb{E}[\sigma'(\mathbf{u}^\top \mathbf{x}) \sigma'(\mathbf{u}^\top \mathbf{z})],\end{aligned}$$

where we used that \tilde{v}_i^2 is independent of $\sigma'(x_i^{(1)}) \sigma'(z_i^{(1)})$. By the same argument

$$\frac{1}{n} \sum_{i=1}^n \sigma(x_i^{(1)}) \sigma(z_i^{(1)}) \rightarrow \mathbb{E}[\sigma(\mathbf{u}^\top \mathbf{x}) \sigma(\mathbf{u}^\top \mathbf{z})]$$

almost surely as $n \rightarrow \infty$. This shows the first statement.

The existence of n_0 follows similarly by an application of Theorem A.22. \square

Example 11.17 (K^{LC} for ReLU). Let $\sigma(x) = \max\{0, x\}$ and let $\mathcal{D}(0, 1)$ be the standard normal distribution. For $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$ denote by

$$\theta = \arccos \left(\frac{\mathbf{x}^\top \mathbf{z}}{\|\mathbf{x}\| \|\mathbf{z}\|} \right)$$

the angle between these vectors. Then according to [36, Appendix A], it holds with $u_i \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1)$, $i = 1, \dots, d$,

$$K^{\text{LC}}(\mathbf{x}, \mathbf{z}) = \mathbb{E}[\sigma(\mathbf{u}^\top \mathbf{x}) \sigma(\mathbf{u}^\top \mathbf{z})] = \frac{\|\mathbf{x}\| \|\mathbf{z}\|}{2\pi d} (\sin(\theta) + (\pi - \theta) \cos(\theta)).$$

11.5.3 Gradient descent

We now proceed similar as in [129, App. G], to show that Theorem 11.13 is applicable to the wide neural network (11.5.1) with high probability under random initialization (11.5.3). This will imply that gradient descent can find global minimizers when training wide neural networks. We work under the following assumptions on the activation function and training data.

Assumption 11.18. There exist $R < \infty$ and $0 < \lambda_{\min}^{\text{LC}} \leq \lambda_{\max}^{\text{LC}} < \infty$ such that

- (a) for the activation function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ holds $|\sigma(0)|, \text{Lip}(\sigma), \text{Lip}(\sigma') \leq R$,
- (b) $\|\mathbf{x}_i\|, |y_i| \leq R$ for all training data $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$, $i = 1, \dots, m$,
- (c) the kernel matrix of the neural tangent kernel

$$(K^{\text{LC}}(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^m \in \mathbb{R}^{m \times m}$$

is regular and its eigenvalues belong to $[\lambda_{\min}^{\text{LC}}, \lambda_{\max}^{\text{LC}}]$.

We start by showing Assumption 11.12 (a) for the present setting. More precisely, we give bounds for the eigenvalues of the empirical tangent kernel.

Lemma 11.19. *Let Assumption 11.18 be satisfied. Then for every $\delta > 0$ there exists $n_0(\delta, \lambda_{\min}^{\text{LC}}, m, R) \in \mathbb{R}$ such that for all $n \geq n_0$ with probability at least $1 - \delta$ all eigenvalues of*

$$(\hat{K}_n(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^m = (\langle \nabla_{\mathbf{w}} \Phi(\mathbf{x}_i, \mathbf{w}_0), \nabla_{\mathbf{w}} \Phi(\mathbf{x}_j, \mathbf{w}_0) \rangle)_{i,j=1}^m \in \mathbb{R}^{m \times m}$$

belong to $[n\lambda_{\min}^{\text{LC}}/2, 2n\lambda_{\max}^{\text{LC}}]$.

Proof. Denote $\hat{\mathbf{G}}_n := (\hat{K}_n(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^m$ and $\mathbf{G}^{\text{LC}} := (K^{\text{LC}}(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^m$. By Theorem 11.16, there exists n_0 such that for all $n \geq n_0$ holds with probability at least $1 - \delta$ that

$$\left\| \mathbf{G}^{\text{LC}} - \frac{1}{n} \hat{\mathbf{G}}_n \right\| \leq \frac{\lambda_{\min}^{\text{LC}}}{2}.$$

Assuming this bound to hold

$$\frac{1}{n} \|\hat{\mathbf{G}}_n\| = \sup_{\substack{\mathbf{a} \in \mathbb{R}^m \\ \|\mathbf{a}\|=1}} \frac{1}{n} \|\hat{\mathbf{G}}_n \mathbf{a}\| \geq \inf_{\substack{\mathbf{a} \in \mathbb{R}^m \\ \|\mathbf{a}\|=1}} \|\mathbf{G}^{\text{LC}} \mathbf{a}\| - \frac{\lambda_{\min}^{\text{LC}}}{2} \geq \lambda_{\min}^{\text{LC}} - \frac{\lambda_{\min}^{\text{LC}}}{2} \geq \frac{\lambda_{\min}^{\text{LC}}}{2},$$

where we have used that $\lambda_{\min}^{\text{LC}}$ is the smallest eigenvalue, and thus singular value, of the symmetric positive definite matrix \mathbf{G}^{LC} . This shows that the smallest eigenvalue of $\hat{\mathbf{G}}_n$ is larger or equal to $\lambda_{\min}^{\text{LC}}/2$. Similarly, we conclude that the largest eigenvalue is bounded from above by $\lambda_{\max}^{\text{LC}} + \lambda_{\min}^{\text{LC}}/2 \leq \lambda_{\max}^{\text{LC}}$. This concludes the proof. \square

Next we check Assumption 11.12 (b). To this end we first bound the norm of a random matrix.

Lemma 11.20. *Let $\mathcal{D}(0, 1)$ be as in Assumption 11.14, and let $\mathbf{W} \in \mathbb{R}^{n \times d}$ with $W_{ij} \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1)$. Denote the fourth moment of $\mathcal{D}(0, 1)$ by μ_4 . Then*

$$\mathbb{P}\left[\|\mathbf{W}\| \leq \sqrt{n(d+1)}\right] \geq 1 - \frac{d\mu_4}{n}.$$

Proof. It holds

$$\|\mathbf{W}\| \leq \|\mathbf{W}\|_F = \left(\sum_{i=1}^n \sum_{j=1}^d W_{ij}^2 \right)^{1/2}.$$

The $\alpha_i := \sum_{j=1}^d W_{ij}^2$, $i = 1, \dots, n$, are i.i.d. distributed with expectation d and finite variance dC , where $C \leq \mu_4$ is the variance of W_{11}^2 . By Theorem A.22

$$\mathbb{P}\left[\|\mathbf{W}\| > \sqrt{n(d+1)}\right] \leq \mathbb{P}\left[\frac{1}{n} \sum_{i=1}^n \alpha_i > d+1\right] \leq \mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n \alpha_i - d\right| > 1\right] \leq \frac{d\mu_4}{n},$$

which concludes the proof. \square

Lemma 11.21. *Let Assumption 11.18 (a) be satisfied with some constant R . Then there exists $M(R)$, and for all $c, \delta > 0$ there exists $n_0(c, d, \delta, R) \in \mathbb{N}$ such that for all $n \geq n_0$ it holds with probability at least $1 - \delta$*

$$\begin{aligned} \|\nabla_{\mathbf{w}}\Phi(\mathbf{x}, \mathbf{w})\| &\leq M\sqrt{n} && \text{for all } \mathbf{w} \in B_{cn^{-1/2}}(\mathbf{w}_0) \\ \|\nabla_{\mathbf{w}}\Phi(\mathbf{x}, \mathbf{w}) - \nabla_{\mathbf{w}}\Phi(\mathbf{x}, \mathbf{v})\| &\leq M\sqrt{n}\|\mathbf{w} - \mathbf{v}\| && \text{for all } \mathbf{w}, \mathbf{v} \in B_{cn^{-1/2}}(\mathbf{w}_0) \end{aligned}$$

for all $\mathbf{x} \in \mathbb{R}^d$ with $\|\mathbf{x}\| \leq R$.

Proof. Due to the initialization (11.5.3), by Lemma 11.20 we can find $n_0(\delta, d)$ such that for all $n \geq n_0$ holds with probability at least $1 - \delta$ that

$$\|\mathbf{v}_0\| \leq 2 \quad \text{and} \quad \|\mathbf{U}_0\| \leq 2\sqrt{n}. \quad (11.5.5)$$

For the rest of this proof we fix arbitrary $\mathbf{x} \in \mathbb{R}^d$ and $n \geq n_0 \geq c^2$ such that

$$\|\mathbf{x}\| \leq R \quad \text{and} \quad n^{-1/2}c \leq 1.$$

We need to show that the claimed inequalities hold as long as (11.5.5) is satisfied. We will several times use that for all $\mathbf{p}, \mathbf{q} \in \mathbb{R}^n$

$$\|\mathbf{p} \odot \mathbf{q}\| \leq \|\mathbf{p}\| \|\mathbf{q}\| \quad \text{and} \quad \|\sigma(\mathbf{p})\| \leq R\sqrt{n} + R\|\mathbf{p}\|$$

since $|\sigma(x)| \leq R \cdot (1 + |x|)$. The same holds for σ' .

Step 1. We show the bound on the gradient. Fix

$$\mathbf{w} = (\mathbf{U}, \mathbf{b}, \mathbf{v}, c) \quad \text{s.t.} \quad \|\mathbf{w} - \mathbf{w}_0\| \leq cn^{-1/2}.$$

Using formula (11.5.2) for $\nabla_{\mathbf{b}}\Phi$ and the above inequalities

$$\begin{aligned} \|\nabla_{\mathbf{b}}\Phi(\mathbf{x}, \mathbf{w})\| &\leq \|\nabla_{\mathbf{b}}\Phi(\mathbf{x}, \mathbf{w}_0)\| + \|\nabla_{\mathbf{b}}\Phi(\mathbf{x}, \mathbf{w}) - \nabla_{\mathbf{b}}\Phi(\mathbf{x}, \mathbf{w}_0)\| \\ &= \|\mathbf{v}_0 \odot \sigma'(\mathbf{U}_0\mathbf{x})\| + \|\mathbf{v} \odot \sigma'(\mathbf{U}\mathbf{x} + \mathbf{b}) - \mathbf{v}_0 \odot \sigma'(\mathbf{U}_0\mathbf{x})\| \\ &\leq 2(R\sqrt{n} + 2R^2\sqrt{n}) + \|\mathbf{v} \odot \sigma'(\mathbf{U}\mathbf{x} + \mathbf{b}) - \mathbf{v}_0 \odot \sigma'(\mathbf{U}_0\mathbf{x})\|. \end{aligned} \quad (11.5.6)$$

Due to

$$\|\mathbf{U}\| \leq \|\mathbf{U}_0\| + \|\mathbf{U}_0 - \mathbf{U}\|_F \leq 2\sqrt{n} + cn^{-1/2} \leq 3\sqrt{n}, \quad (11.5.7)$$

the last norm in (11.5.6) is bounded by

$$\begin{aligned} &\|(\mathbf{v} - \mathbf{v}_0) \odot \sigma'(\mathbf{U}\mathbf{x} + \mathbf{b})\| + \|\mathbf{v}_0 \odot (\sigma'(\mathbf{U}\mathbf{x} + \mathbf{b}) - \sigma'(\mathbf{U}_0\mathbf{x}))\| \\ &\leq cn^{-1/2}(R\sqrt{n} + R \cdot (\|\mathbf{U}\|\|\mathbf{x}\| + \|\mathbf{b}\|)) + 2R \cdot (\|\mathbf{U} - \mathbf{U}_0\|\|\mathbf{x}\| + \|\mathbf{b}\|) \\ &\leq R\sqrt{n} + 3\sqrt{n}R^2 + cn^{-1/2}R + 2R \cdot (cn^{-1/2}R + cn^{-1/2}) \\ &\leq \sqrt{n}(4R + 5R^2) \end{aligned}$$

and therefore

$$\|\nabla_{\mathbf{b}}\Phi(\mathbf{x}, \mathbf{w})\| \leq \sqrt{n}(6R + 9R^2).$$

For the gradient with respect to \mathbf{U} we use $\nabla_{\mathbf{U}}\Phi(\mathbf{x}, \mathbf{w}) = \nabla_{\mathbf{b}}\Phi(\mathbf{x}, \mathbf{w})\mathbf{x}^\top$, so that

$$\|\nabla_{\mathbf{U}}\Phi(\mathbf{x}, \mathbf{w})\|_F = \|\nabla_{\mathbf{b}}\Phi(\mathbf{x}, \mathbf{w})\mathbf{x}^\top\|_F = \|\nabla_{\mathbf{b}}\Phi(\mathbf{x}, \mathbf{w})\|\|\mathbf{x}\| \leq \sqrt{n}(6R^2 + 9R^3).$$

Next

$$\begin{aligned} \|\nabla_{\mathbf{v}}\Phi(\mathbf{x}, \mathbf{w})\| &= \|\sigma(\mathbf{U}\mathbf{x} + \mathbf{b})\| \\ &\leq R\sqrt{n} + R\|\mathbf{U}\mathbf{x} + \mathbf{b}\| \\ &\leq R\sqrt{n} + R \cdot (3\sqrt{n}R + cn^{-1/2}) \\ &\leq \sqrt{n}(2R + 3R^2), \end{aligned}$$

and finally $\nabla_c\Phi(\mathbf{x}, \mathbf{w}) = 1$. In all, with $M_1(R) := (1 + 8R + 12R^2)$

$$\|\nabla_{\mathbf{w}}\Phi(\mathbf{x}, \tilde{\mathbf{w}})\| \leq \sqrt{n}M_1(R).$$

Step 2. We show Lipschitz continuity. Fix

$$\mathbf{w} = (\mathbf{U}, \mathbf{b}, \mathbf{v}, c) \quad \text{and} \quad \tilde{\mathbf{w}} = (\tilde{\mathbf{U}}, \tilde{\mathbf{b}}, \tilde{\mathbf{v}}, \tilde{c})$$

such that $\|\mathbf{w} - \mathbf{w}_0\|, \|\tilde{\mathbf{w}} - \mathbf{w}_0\| \leq cn^{-1/2}$. Then

$$\|\nabla_{\mathbf{b}}\Phi(\mathbf{x}, \mathbf{w}) - \nabla_{\mathbf{b}}\Phi(\mathbf{x}, \tilde{\mathbf{w}})\| = \|\mathbf{v} \odot \sigma'(\mathbf{U}\mathbf{x} + \mathbf{b}) - \tilde{\mathbf{v}} \odot \sigma'(\tilde{\mathbf{U}}\mathbf{x} + \tilde{\mathbf{b}})\|.$$

Using $\|\tilde{\mathbf{v}}\| \leq \|\mathbf{v}_0\| + cn^{-1/2} \leq 3$ and (11.5.7), this term is bounded by

$$\begin{aligned} &\|(\mathbf{v} - \tilde{\mathbf{v}}) \odot \sigma'(\mathbf{U}\mathbf{x} + \mathbf{b})\| + \|\tilde{\mathbf{v}} \odot (\sigma'(\mathbf{U}\mathbf{x} + \mathbf{b}) - \sigma'(\tilde{\mathbf{U}}\mathbf{x} + \tilde{\mathbf{b}}))\| \\ &\leq \|\mathbf{v} - \tilde{\mathbf{v}}\|(R\sqrt{n} + R \cdot (\|\mathbf{U}\|\|\mathbf{x}\| + \|\mathbf{b}\|)) + 3R \cdot (\|\mathbf{x}\|\|\mathbf{U} - \tilde{\mathbf{U}}\| + \|\mathbf{b} - \tilde{\mathbf{b}}\|) \\ &\leq \|\mathbf{w} - \tilde{\mathbf{w}}\|\sqrt{n}(5R + 6R^2). \end{aligned}$$

For $\nabla_{\mathbf{U}}\Phi(\mathbf{x}, \mathbf{w})$ we obtain similar as in Step 1

$$\begin{aligned} \|\nabla_{\mathbf{U}}\Phi(\mathbf{x}, \mathbf{w}) - \nabla_{\mathbf{U}}\Phi(\mathbf{x}, \tilde{\mathbf{w}})\|_F &= \|\mathbf{x}\|\|\nabla_{\mathbf{b}}\Phi(\mathbf{x}, \mathbf{w}) - \nabla_{\mathbf{b}}\Phi(\mathbf{x}, \tilde{\mathbf{w}})\| \\ &\leq \|\mathbf{w} - \tilde{\mathbf{w}}\|\sqrt{n}(5R^2 + 6R^3). \end{aligned}$$

Next

$$\begin{aligned} \|\nabla_{\mathbf{v}}\Phi(\mathbf{x}, \mathbf{w}) - \nabla_{\mathbf{v}}\Phi(\mathbf{x}, \tilde{\mathbf{w}})\| &= \|\sigma(\mathbf{U}\mathbf{x} + \mathbf{b}) - \sigma(\tilde{\mathbf{U}}\mathbf{x} + \tilde{\mathbf{b}})\| \\ &\leq R \cdot (\|\mathbf{U} - \tilde{\mathbf{U}}\|\|\mathbf{x}\| + \|\mathbf{b} - \tilde{\mathbf{b}}\|) \\ &\leq \|\mathbf{w} - \tilde{\mathbf{w}}\|(R^2 + R) \end{aligned}$$

and finally $\nabla_c\Phi(\mathbf{x}, \mathbf{w}) = 1$ is constant. With $M_2(R) := R + 6R^2 + 6R^3$ this shows

$$\|\nabla_{\mathbf{w}}\Phi(\mathbf{x}, \mathbf{w}) - \nabla_{\mathbf{w}}\Phi(\mathbf{x}, \tilde{\mathbf{w}})\| \leq \sqrt{n}M_2(R)\|\mathbf{w} - \tilde{\mathbf{w}}\|.$$

In all, this concludes the proof with $M(R) := \max\{M_1(R), M_2(R)\}$. \square

Before coming to the main result of this section, we first show that the initial error $f(\mathbf{w}_0)$ remains bounded with high probability.

Lemma 11.22. *Let Assumption 11.18 (a), (b) be satisfied. Then for every $\delta > 0$ exists $R_0(\delta, m, R) > 0$ such that for all $n \in \mathbb{N}$*

$$\mathbb{P}[f(\mathbf{w}_0) \leq R_0] \geq 1 - \delta.$$

Proof. Let $i \in \{1, \dots, m\}$, and set $\boldsymbol{\alpha} := \mathbf{U}_0 \mathbf{x}_i$ and $\tilde{v}_j := \sqrt{n} v_{0,j}$ for $j = 1, \dots, n$, so that $\tilde{v}_j \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1)$. Then

$$\Phi(\mathbf{x}_i, \mathbf{w}_0) = \frac{1}{\sqrt{n}} \sum_{j=1}^n \tilde{v}_j \sigma(\alpha_j).$$

By Assumption 11.14 and (11.5.3), the $\tilde{v}_j \sigma(\alpha_j)$, $j = 1, \dots, n$, are i.i.d. centered random variables with finite variance bounded by a constant $C(R)$ independent of n . Thus the variance of $\Phi(\mathbf{x}_i, \mathbf{w}_0)$ is also bounded by $C(R)$. By Chebyshev's inequality, see Lemma A.21, for every $k > 0$

$$\mathbb{P}[|\Phi(\mathbf{x}_i, \mathbf{w}_0)| \geq k\sqrt{C}] \leq \frac{1}{k^2}.$$

Setting $k = \sqrt{m/\delta}$

$$\begin{aligned} \mathbb{P}\left[\sum_{i=1}^m |\Phi(\mathbf{x}_i, \mathbf{w}_0) - y_i|^2 \geq m(k\sqrt{C} + R)^2\right] &\leq \sum_{i=1}^m \mathbb{P}\left[|\Phi(\mathbf{x}_i, \mathbf{w}_0) - y_i| \geq k\sqrt{C} + R\right] \\ &\leq \sum_{i=1}^m \mathbb{P}\left[|\Phi(\mathbf{x}_i, \mathbf{w}_0)| \geq k\sqrt{C}\right] \leq \delta, \end{aligned}$$

which shows the claim with $R_0 = m \cdot (\sqrt{Cm/\delta} + R)^2$. \square

The next theorem is the main result of this section. It states that in the present setting gradient descent converges to a global minimizer and the limiting network achieves zero loss, i.e. interpolates the data. Moreover, during training the network weights remain close to initialization if the network width n is large.

Theorem 11.23. *Let Assumption 11.18 be satisfied, and let the parameters \mathbf{w}_0 of the neural network Φ in (11.5.1) be initialized according to (11.5.3). Fix a learning rate*

$$h < \frac{2}{\lambda_{\min}^{\text{LC}} + 4\lambda_{\max}^{\text{LC}}} \frac{1}{n}$$

and with the objective function (11.0.1b) let for all $k \in \mathbb{N}$

$$\mathbf{w}_{k+1} = \mathbf{w}_k - h \nabla f(\mathbf{w}_k).$$

Then for every $\delta > 0$ there exist $C > 0$, $n_0 \in \mathbb{N}$ such that for all $n \geq n_0$ holds with probability at least $1 - \delta$ that for all $k \in \mathbb{N}$

$$\begin{aligned}\|\mathbf{w}_k - \mathbf{w}_0\| &\leq \frac{C}{\sqrt{n}} \\ f(\mathbf{w}_k) &\leq C \left(1 - \frac{hn}{2\lambda_{\min}^{\text{LC}}}\right)^{2k}.\end{aligned}$$

Proof. We wish to apply Theorem 11.13, which requires Assumption 11.12 to be satisfied. By Lemma 11.19, 11.21 and 11.22, for every $c > 0$ we can find n_0 such that for all $n \geq n_0$ with probability at least $1 - \delta$ we have $\sqrt{f(\mathbf{w}_0)} \leq \sqrt{R_0}$ and Assumption 11.12 (a), (b) holds with the values

$$L = M\sqrt{n}, \quad U = M\sqrt{n}, \quad r = cn^{-1/2}, \quad \lambda_{\min} = \frac{n\lambda_{\min}^{\text{LC}}}{2}, \quad \lambda_{\max} = 2n\lambda_{\max}^{\text{LC}}.$$

For Assumption 11.12 (c), it suffices that

$$M\sqrt{n} \leq \frac{n^2(\lambda_{\min}^{\text{LC}}/2)^2}{12m^{3/2}M^2n\sqrt{R_0}} \quad \text{and} \quad cn^{-1/2} \geq \frac{2mM\sqrt{n}}{n}\sqrt{R_0}.$$

Choosing $c > 0$ and n large enough, the inequalities hold. The statement is now a direct consequence of Theorem 11.13. \square

11.5.4 Proximity to linearized model

The analysis thus far was based on the linearization Φ^{lin} describing the behaviour of the full network Φ well in a neighbourhood of the initial parameters \mathbf{w}_0 . Moreover, Theorem 11.23 states that the parameters remain in an $O(n^{-1/2})$ neighbourhood of \mathbf{w}_0 during training. This suggests that the trained full model $\lim_{k \rightarrow \infty} \Phi(\mathbf{x}, \mathbf{w}_k)$ yields predictions similar to the trained linearized model.

To describe this phenomenon, we adopt again the notations $\Phi^{\text{lin}} : \mathbb{R}^d \times \mathbb{R}^n \rightarrow \mathbb{R}$ and f^{lin} from (11.3.1) and (11.3.3). Initializing \mathbf{w}_0 according to (11.5.3) and setting $\mathbf{p}_0 = \mathbf{w}_0$, gradient descent computes the parameter updates

$$\mathbf{w}_{k+1} = \mathbf{w}_k - h\nabla_{\mathbf{w}}f(\mathbf{w}_k), \quad \mathbf{p}_{k+1} = \mathbf{p}_k - h\nabla_{\mathbf{w}}f^{\text{lin}}(\mathbf{p}_k)$$

for the full and linearized models, respectively. Let us consider the dynamics of the prediction of the network on the training data. Writing

$$\Phi(\mathbf{X}, \mathbf{w}) := (\Phi(\mathbf{x}_i, \mathbf{w}))_{i=1}^m \in \mathbb{R}^m \quad \text{such that} \quad \nabla_{\mathbf{w}}\Phi(\mathbf{X}, \mathbf{w}) \in \mathbb{R}^{m \times n}$$

it holds

$$\nabla_{\mathbf{w}}f(\mathbf{w}) = \nabla_{\mathbf{w}}\|\Phi(\mathbf{X}, \mathbf{w}) - \mathbf{y}\|^2 = 2\nabla_{\mathbf{w}}\Phi(\mathbf{X}, \mathbf{w})^\top (\Phi(\mathbf{X}, \mathbf{w}) - \mathbf{y}).$$

Thus for the full model

$$\begin{aligned}\Phi(\mathbf{X}, \mathbf{w}_{k+1}) &= \Phi(\mathbf{X}, \mathbf{w}_k) + \nabla_{\mathbf{w}}\Phi(\mathbf{X}, \tilde{\mathbf{w}}_k)(\mathbf{w}_{k+1} - \mathbf{w}_k) \\ &= \Phi(\mathbf{X}, \mathbf{w}_k) - 2h\nabla_{\mathbf{w}}\Phi(\mathbf{X}, \tilde{\mathbf{w}}_k)\nabla_{\mathbf{w}}\Phi(\mathbf{X}, \mathbf{w}_k)^\top (\Phi(\mathbf{X}, \mathbf{w}_k) - \mathbf{y}),\end{aligned}\tag{11.5.8}$$

where $\tilde{\mathbf{w}}_k$ is in the convex hull of \mathbf{w}_k and \mathbf{w}_{k+1} .

Similarly, for the linearized model with (cp. (11.3.1))

$$\Phi^{\text{lin}}(\mathbf{X}, \mathbf{w}) := (\Phi^{\text{lin}}(\mathbf{x}_i, \mathbf{w}))_{i=1}^m \in \mathbb{R}^m \quad \text{and} \quad \nabla_{\mathbf{p}} \Phi^{\text{lin}}(\mathbf{X}, \mathbf{p}) = \nabla_{\mathbf{w}} \Phi(\mathbf{X}, \mathbf{w}_0) \in \mathbb{R}^{m \times n}$$

such that

$$\nabla_{\mathbf{p}} f^{\text{lin}}(\mathbf{p}) = \nabla_{\mathbf{p}} \|\Phi^{\text{lin}}(\mathbf{X}, \mathbf{p}) - \mathbf{y}\|^2 = 2 \nabla_{\mathbf{w}} \Phi(\mathbf{X}, \mathbf{w}_0)^\top (\Phi^{\text{lin}}(\mathbf{X}, \mathbf{p}) - \mathbf{y})$$

and

$$\begin{aligned} \Phi^{\text{lin}}(\mathbf{X}, \mathbf{p}_{k+1}) &= \Phi^{\text{lin}}(\mathbf{X}, \mathbf{p}_k) + \nabla_{\mathbf{p}} \Phi^{\text{lin}}(\mathbf{X}, \mathbf{p}_0)(\mathbf{p}_{k+1} - \mathbf{p}_k) \\ &= \Phi^{\text{lin}}(\mathbf{X}, \mathbf{p}_k) - 2h \nabla_{\mathbf{w}} \Phi(\mathbf{X}, \mathbf{w}_0) \nabla_{\mathbf{w}} \Phi(\mathbf{X}, \mathbf{w}_0)^\top (\Phi^{\text{lin}}(\mathbf{X}, \mathbf{p}_k) - \mathbf{y}). \end{aligned} \quad (11.5.9)$$

Remark 11.24. From (11.5.9) it is easy to see that with $\mathbf{A} := 2h \nabla_{\mathbf{w}} \Phi(\mathbf{X}, \mathbf{w}_0) \nabla_{\mathbf{w}} \Phi(\mathbf{X}, \mathbf{w}_0)^\top$ and $\mathbf{B} := \mathbf{I}_m - \mathbf{A}$ holds the explicit formula

$$\Phi^{\text{lin}}(\mathbf{X}, \mathbf{p}_k) = \mathbf{B}^k \Phi^{\text{lin}}(\mathbf{X}, \mathbf{p}_0) + \sum_{j=0}^{k-1} \mathbf{B}^j \mathbf{A} \mathbf{y}$$

for the prediction of the linear model in step k . Note that if \mathbf{A} is regular and h is small enough, then \mathbf{B}^k converges to the zero matrix as $k \rightarrow \infty$ and $\sum_{j=0}^{\infty} \mathbf{B}^j = \mathbf{A}^{-1}$ since this is a Neumann series.

Comparing the two dynamics (11.5.8) and (11.5.9), the difference only lies in the two $\mathbb{R}^{m \times m}$ matrices

$$2h \nabla_{\mathbf{w}} \Phi(\mathbf{X}, \tilde{\mathbf{w}}_k) \nabla_{\mathbf{w}} \Phi(\mathbf{X}, \mathbf{w}_k)^\top \quad \text{and} \quad 2h \nabla_{\mathbf{w}} \Phi(\mathbf{X}, \mathbf{w}_0) \nabla_{\mathbf{w}} \Phi(\mathbf{X}, \mathbf{w}_0)^\top.$$

Recall that the step size h in Theorem 11.23 scales like $1/n$.

Proposition 11.25. *Consider the setting of Theorem 11.23. Then there exists $C < \infty$, and for every $\delta > 0$ there exists n_0 such that for all $n \geq n_0$ holds with probability at least $1 - \delta$ that for all $k \in \mathbb{N}$*

$$\frac{1}{n} \|\nabla_{\mathbf{w}} \Phi(\mathbf{X}, \tilde{\mathbf{w}}_k) \nabla_{\mathbf{w}} \Phi(\mathbf{X}, \mathbf{w}_k)^\top - \nabla_{\mathbf{p}} \Phi(\mathbf{X}, \mathbf{p}_0) \nabla_{\mathbf{p}} \Phi(\mathbf{X}, \mathbf{p}_0)^\top\| \leq C n^{-1/2}.$$

Proof. Consider the setting of the proof of Theorem 11.23. Then for every $k \in \mathbb{N}$ holds $\|\mathbf{w}_k - \mathbf{w}_0\| \leq r$ and thus also $\|\tilde{\mathbf{w}}_k - \mathbf{w}_0\| \leq r$, where $r = c n^{-1/2}$. Thus Lemma 11.21 implies the norm to be bounded by

$$\begin{aligned} &\frac{1}{n} \|\nabla_{\mathbf{w}} \Phi(\mathbf{X}, \tilde{\mathbf{w}}_k) - \nabla_{\mathbf{p}} \Phi(\mathbf{X}, \mathbf{p}_0)\| \|\nabla_{\mathbf{w}} \Phi(\mathbf{X}, \mathbf{w}_k)^\top\| + \\ &\quad \frac{1}{n} \|\nabla_{\mathbf{p}} \Phi(\mathbf{X}, \mathbf{p}_0)\| \|\nabla_{\mathbf{w}} \Phi(\mathbf{X}, \mathbf{w}_k)^\top - \nabla_{\mathbf{p}} \Phi(\mathbf{X}, \mathbf{p}_0)^\top\| \\ &\leq m M (\|\tilde{\mathbf{w}}_k - \mathbf{p}_0\| + \|\mathbf{w}_k - \mathbf{p}_0\|) \leq c m M n^{-1/2} \end{aligned}$$

which gives the statement. \square

By Proposition 11.25 the two matrices driving the dynamics (11.5.8) and (11.5.9) remain in an $O(n^{-1/2})$ neighbourhood of each other throughout training. This allows to show the following proposition, which states that the prediction function learned by the network gets arbitrarily close to the one learned by the linearized version in the limit $n \rightarrow \infty$. The proof, which we omit, is based on Grönwall's inequality. See [105, 129].

Proposition 11.26. *Consider the setting of Theorem 11.23. Then there exists $C < \infty$, and for every $\delta > 0$ there exists n_0 such that for all $n \geq n_0$ holds with probability at least $1 - \delta$ that for all $\|\mathbf{x}\| \leq 1$*

$$\sup_{k \in \mathbb{N}} |\Phi(\mathbf{x}, \mathbf{w}_k) - \Phi^{\text{lin}}(\mathbf{x}, \mathbf{p}_k)| \leq Cn^{-1/2}.$$

11.5.5 Connection to Gaussian processes

In the previous section, we established that for large widths, the trained neural network mirrors the behaviour of the trained linearized model, which itself is closely connected to kernel least-squares with the neural tangent kernel. Yet, as pointed out in Remark 11.4, the obtained model still strongly depends on the choice of random initialization $\mathbf{w}_0 \in \mathbb{R}^n$. We should thus understand both the model at initialization $\mathbf{x} \mapsto \Phi(\mathbf{x}, \mathbf{w}_0)$ and the model after training $\mathbf{x} \mapsto \Phi(\mathbf{x}, \mathbf{w}_k)$, as random draws of a certain distribution over functions. To make this precise, let us introduce Gaussian processes.

Definition 11.27. Let (Ω, \mathbb{P}) be a probability space, and let $g : \mathbb{R}^d \times \Omega \rightarrow \mathbb{R}$. We call g a **Gaussian process** with mean function $m : \mathbb{R}^d \rightarrow \mathbb{R}$ and covariance function $c : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ if

- (a) for each $\mathbf{x} \in \mathbb{R}^d$ holds $\omega \mapsto g(\mathbf{x}, \omega)$ is a random variable,
- (b) for all $k \in \mathbb{N}$ and all $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^d$ the random variables $g(\mathbf{x}_1, \cdot), \dots, g(\mathbf{x}_k, \cdot)$ have a joint Gaussian distribution such that

$$(g(\mathbf{x}_1, \omega), \dots, g(\mathbf{x}_k, \omega)) \sim \mathcal{N}\left(m(\mathbf{x}_i)_{i=1}^k, (c(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^k\right).$$

In words, g is a Gaussian process, if $\omega \mapsto g(\mathbf{x}, \omega)$ defines a collection of random variables indexed over $\mathbf{x} \in \mathbb{R}^d$, such that the joint distribution of $(g(\mathbf{x}_1, \cdot))_{j=1}^n$ is a Gaussian whose mean and variance are determined by m and c respectively. Fixing $\omega \in \Omega$, we can then interpret $\mathbf{x} \mapsto g(\mathbf{x}, \omega)$ as a random draw from a distribution over functions.

As first observed in [155], certain neural networks at initialization tend to Gaussian processes in the infinite width limit.

Proposition 11.28. Consider depth- n networks Φ_n as in (11.5.1) with initialization (11.5.3), and define with $u_i \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1/d)$, $i = 1, \dots, d$,

$$c(\mathbf{x}, \mathbf{z}) := \mathbb{E}[\sigma(\mathbf{u}^\top \mathbf{x})\sigma(\mathbf{u}^\top \mathbf{z})] \quad \text{for all } \mathbf{x}, \mathbf{z} \in \mathbb{R}^d.$$

Then for all distinct $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^d$ it holds that

$$\lim_{n \rightarrow \infty} (\Phi_n(\mathbf{x}_1, \mathbf{w}_0), \dots, \Phi_n(\mathbf{x}_k, \mathbf{w}_0)) \sim \mathbf{N}(\mathbf{0}, (c(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^k)$$

with weak convergence.

Proof. Set $\tilde{v}_i := \sqrt{n}v_{0,i}$ and $\tilde{\mathbf{u}}_i = (U_{0,i1}, \dots, U_{0,id}) \in \mathbb{R}^d$, so that $\tilde{v}_i \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1)$, and the $\tilde{\mathbf{u}}_i \in \mathbb{R}^d$ are also i.i.d., with each component distributed according to $\mathcal{D}(0, 1/d)$.

Then for any $\mathbf{x}_1, \dots, \mathbf{x}_k$

$$\mathbf{Z}_i := \begin{pmatrix} \tilde{v}_i \sigma(\tilde{\mathbf{u}}_i^\top \mathbf{x}_1) \\ \vdots \\ \tilde{v}_i \sigma(\tilde{\mathbf{u}}_i^\top \mathbf{x}_k) \end{pmatrix} \in \mathbb{R}^k \quad i = 1, \dots, n,$$

defines n centered i.i.d. vectors in \mathbb{R}^k . By the central limit theorem, see Theorem A.24,

$$\begin{pmatrix} \Phi(\mathbf{x}_1, \mathbf{w}_0) \\ \vdots \\ \Phi(\mathbf{x}_k, \mathbf{w}_0) \end{pmatrix} = \frac{1}{\sqrt{n}} \sum_{j=1}^n \mathbf{Z}_j$$

converges weakly to $\mathbf{N}(\mathbf{0}, \mathbf{C})$, where

$$C_{ij} = \mathbb{E}[\tilde{v}_i^2 \sigma(\tilde{\mathbf{u}}_i^\top \mathbf{x}_i) \sigma(\tilde{\mathbf{u}}_j^\top \mathbf{x}_j)] = \mathbb{E}[\sigma(\tilde{\mathbf{u}}_i^\top \mathbf{x}_i) \sigma(\tilde{\mathbf{u}}_j^\top \mathbf{x}_j)].$$

This concludes the proof. \square

In the sense of Proposition 11.28, the network $\Phi(\mathbf{x}, \mathbf{w}_0)$ converges to a Gaussian process as the width n tends to infinity. Using the explicit dynamics of the linearized network outlined in Remark 11.24, one can show that the linearized network after training also corresponds to a Gaussian process (for some mean and covariance function depending on the data, the architecture, and the initialization). As the full and linearized models converge in the infinite width limit, we can infer that wide networks post-training resemble draws from a Gaussian process, see [129, Sec. 2.3.1] and [45].

Rather than delving into the technical details of such statements, in Figure 11.3 we plot 80 different realizations of a neural network before and after training, i.e.

$$\mathbf{x} \mapsto \Phi(\mathbf{x}, \mathbf{w}_0) \quad \text{and} \quad \mathbf{x} \mapsto \Phi(\mathbf{x}, \mathbf{w}_k). \quad (11.5.10)$$

We chose the architecture as (11.5.1) with activation function $\sigma = \arctan(x)$, width $n = 250$ and initialization

$$U_{0;ij} \stackrel{\text{iid}}{\sim} \mathbf{N}\left(0, \frac{3}{d}\right), \quad v_{0;i} \stackrel{\text{iid}}{\sim} \mathbf{N}\left(0, \frac{3}{n}\right), \quad b_{0;i}, c_0 \stackrel{\text{iid}}{\sim} \mathbf{N}(0, 2). \quad (11.5.11)$$

The network was trained on a dataset of size $m = 3$ with $k = 1000$ steps of gradient descent and constant step size $h = 1/n$. Before training, the network's outputs resemble random draws from a Gaussian process with a constant zero mean function. Post-training, the outputs show minimal variance at the data points, since they essentially interpolate the data, cp. Remark 11.4 and (11.2.4). They exhibit increased variance further from these points, with the precise amount depending on the initialization variance chosen in (11.5.11).

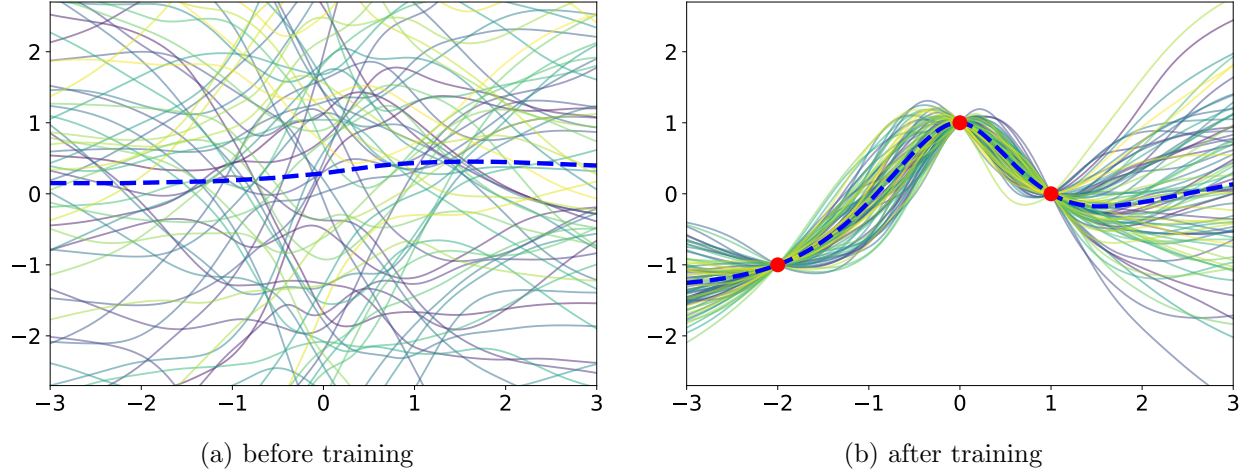


Figure 11.3: 80 realizations of a neural network at initialization (a) and after training on the red data points (b). The blue dashed line shows the mean. Figure based on [129, Fig. 2].

11.6 Normalized initialization

Consider the gradient $\nabla_{\mathbf{w}}\Phi(\mathbf{x}, \mathbf{w}_0)$ as in (11.5.2) with LeCun initialization. Since the components of \mathbf{v} behave like $v_i \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1/n)$, it is easy to check that in terms of the width n

$$\begin{aligned} \mathbb{E}[\|\nabla_{\mathbf{U}}\Phi(\mathbf{x}, \mathbf{w}_0)\|] &= \mathbb{E}[\|(\mathbf{v} \odot \sigma'(\mathbf{U}\mathbf{x} + \mathbf{b}))\mathbf{x}^\top\|] &&= O(1) \\ \mathbb{E}[\|\nabla_{\mathbf{b}}\Phi(\mathbf{x}, \mathbf{w}_0)\|] &= \mathbb{E}[\|\mathbf{v} \odot \sigma'(\mathbf{U}\mathbf{x} + \mathbf{b})\|] &&= O(1) \\ \mathbb{E}[\|\nabla_{\mathbf{v}}\Phi(\mathbf{x}, \mathbf{w}_0)\|] &= \mathbb{E}[\|\sigma(\mathbf{U}\mathbf{x} + \mathbf{b})\|] &&= O(n) \\ \mathbb{E}[\|\nabla_c\Phi(\mathbf{x}, \mathbf{w}_0)\|] &= \mathbb{E}[1] &&= O(1). \end{aligned}$$

As a result of this different scaling, gradient descent with step width $O(n^{-1})$ as in Theorem 11.23, will primarily train the weights \mathbf{v} in the output layer, and will barely move the remaining parameters \mathbf{U} , \mathbf{b} , and c . This is also reflected in the expression for the obtained kernel K^{LC} computed in Theorem 11.16, which corresponds to the contribution of the term $\langle \nabla_{\mathbf{v}}\Phi, \nabla_{\mathbf{v}}\Phi \rangle$.

Remark 11.29. For optimization methods such as ADAM, which scale each component of the gradient individually, the same does not hold in general.

LeCun initialization aims to normalize the variance of the output of all nodes at initialization (the forward dynamics). To also normalize the variance of the gradients (the backward dynamics), in this section we shortly discuss a different architecture and initialization, consistent with the one used in the original NTK paper [105].

11.6.1 Architecture

Let $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}$ be a depth-one neural network

$$\Phi(\mathbf{x}, \mathbf{w}) = \frac{1}{\sqrt{n}} \mathbf{v}^\top \sigma\left(\frac{1}{\sqrt{d}} \mathbf{U} \mathbf{x} + \mathbf{b}\right) + c, \quad (11.6.1)$$

with input $\mathbf{x} \in \mathbb{R}^d$ and parameters $\mathbf{U} \in \mathbb{R}^{n \times d}$, $\mathbf{v} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^n$ and $c \in \mathbb{R}$. We initialize the weights randomly according to $\mathbf{w}_0 = (\mathbf{U}_0, \mathbf{b}_0, \mathbf{v}_0, c_0)$ with parameters

$$U_{0;ij} \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1), \quad v_{0;i} \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1), \quad b_{0;i}, c_0 = 0. \quad (11.6.2)$$

At initialization, (11.6.1), (11.6.2) is equivalent to (11.5.1), (11.5.3). However, for the gradient we obtain

$$\begin{aligned} \nabla_{\mathbf{U}} \Phi(\mathbf{x}, \mathbf{w}) &= n^{-1/2} \left(\mathbf{v} \odot \sigma'(d^{-1/2} \mathbf{U} \mathbf{x} + \mathbf{b}) \right) d^{-1/2} \mathbf{x}^\top \in \mathbb{R}^{n \times d} \\ \nabla_{\mathbf{b}} \Phi(\mathbf{x}, \mathbf{w}) &= n^{-1/2} \mathbf{v} \odot \sigma'(d^{-1/2} \mathbf{U} \mathbf{x} + \mathbf{b}) \in \mathbb{R}^n \\ \nabla_{\mathbf{v}} \Phi(\mathbf{x}, \mathbf{w}) &= n^{-1/2} \sigma(d^{-1/2} \mathbf{U} \mathbf{x} + \mathbf{b}) \in \mathbb{R}^n \\ \nabla_c \Phi(\mathbf{x}, \mathbf{w}) &= 1 \in \mathbb{R}. \end{aligned} \quad (11.6.3)$$

Contrary to (11.5.2), the three gradients with $O(n)$ entries are all scaled by the factor $n^{-1/2}$. This leads to a different training dynamics.

11.6.2 Neural tangent kernel

We compute again the neural tangent kernel. Unlike for LeCun initialization, there is no $1/n$ scaling required to obtain convergence of

$$\hat{K}_n(\mathbf{x}, \mathbf{z}) = \langle \nabla_{\mathbf{w}} \Phi(\mathbf{x}, \mathbf{w}_0), \nabla_{\mathbf{w}} \Phi(\mathbf{z}, \mathbf{w}_0) \rangle$$

as $n \rightarrow \infty$. Here and in the following we consider the setting (11.6.1)–(11.6.2) for Φ and \mathbf{w}_0 . This is also referred to as the NTK initialization, we denote the kernel by K^{NTK} . Due to the different training dynamics, we obtain additional terms in the NTK compared to Theorem 11.23.

Theorem 11.30. *Let $R < \infty$ such that $|\sigma(x)| \leq R \cdot (1 + |x|)$ and $|\sigma'(x)| \leq R \cdot (1 + |x|)$ for all $x \in \mathbb{R}$, and let \mathcal{D} satisfy Assumption 11.14. For any $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$ and $u_i \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1/d)$, $i = 1, \dots, d$, it then holds*

$$\begin{aligned} \lim_{n \rightarrow \infty} \hat{K}_n(\mathbf{x}, \mathbf{z}) &= \left(1 + \frac{\mathbf{x}^\top \mathbf{z}}{d}\right) \mathbb{E}[\sigma'(\mathbf{u}^\top \mathbf{x})^\top \sigma'(\mathbf{u}^\top \mathbf{z})] + \mathbb{E}[\sigma(\mathbf{u}^\top \mathbf{x})^\top \sigma(\mathbf{u}^\top \mathbf{z})] + 1 \\ &=: K^{\text{NTK}}(\mathbf{x}, \mathbf{z}) \end{aligned}$$

almost surely.

Proof. Denote $\mathbf{x}^{(1)} = \mathbf{U}_0 \mathbf{x} + \mathbf{b}_0 \in \mathbb{R}^n$ and $\mathbf{z}^{(1)} = \mathbf{U}_0 \mathbf{z} + \mathbf{b}_0 \in \mathbb{R}^n$. Due to the initialization (11.6.2) and our assumptions on $\mathcal{D}(0, 1)$, the components

$$x_i^{(1)} = \sum_{j=1}^d U_{0;ij} x_j \sim \mathbf{u}^\top \mathbf{x} \quad i = 1, \dots, n$$

are i.i.d. with finite p th moment (independent of n) for all $1 \leq p \leq 8$, and the same holds for the $(\sigma(x_i^{(1)}))_{i=1}^n$, $(\sigma'(x_i^{(1)}))_{i=1}^n$, $(\sigma(z_i^{(1)}))_{i=1}^n$, and $(\sigma'(z_i^{(1)}))_{i=1}^n$.

Then

$$\hat{K}_n(\mathbf{x}, \mathbf{z}) = \left(1 + \frac{\mathbf{x}^\top \mathbf{z}}{d}\right) \frac{1}{n} \sum_{i=1}^n v_i^2 \sigma'(x_i^{(1)}) \sigma'(z_i^{(1)}) + \frac{1}{n} \sum_{i=1}^n \sigma(x_i^{(1)}) \sigma(z_i^{(1)}) + 1.$$

By the law of large numbers and because $\mathbb{E}[v_i^2] = 1$, this converges almost surely to $K^{\text{NTK}}(\mathbf{x}, \mathbf{z})$.

The existence of n_0 follows similarly by an application of Theorem A.22. \square

Example 11.31 (K^{NTK} for ReLU). Let $\sigma(x) = \max\{0, x\}$ and let $\mathcal{D}(0, 1/d)$ be the centered normal distribution on \mathbb{R} with variance $1/d$. For $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$ holds by [36, Appendix A] (also see Exercise 11.36), that with $u_i \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1/d)$, $i = 1, \dots, d$,

$$\mathbb{E}[\sigma'(\mathbf{u}^\top \mathbf{x}) \sigma'(\mathbf{u}^\top \mathbf{z})] = \frac{\pi - \arccos\left(\frac{\mathbf{x}^\top \mathbf{z}}{\|\mathbf{x}\| \|\mathbf{z}\|}\right)}{2\pi}.$$

Together with Example 11.17, this yields an explicit formula for K^{NTK} in Theorem 11.30.

For this network architecture and under suitable assumptions on \mathcal{D} , similar arguments as in Section 11.5 can be used to show convergence of gradient descent to a global minimizer and proximity of the full to the linearized model. We refer to the literature in the bibliography section.

Bibliography and further reading

The discussion on linear and kernel regression in Sections 11.1 and 11.2 is quite standard, and can similarly be found in many textbooks. For more details on kernel methods we refer for instance to [41, 205]. The neural tangent kernel and its connection to the training dynamics was first investigated in [105] using an architecture similar to the one in Section 11.6. Since then, many works have extended this idea and presented differing perspectives on the topic, see for instance [2, 55, 5, 35]. Our presentation in Sections 11.4, 11.5, and 11.6 primarily follows [129] who also discussed the case of LeCun initialization. Especially for the main results in Theorem 11.13 and Theorem 11.23, we largely follow the arguments in this paper. The above references additionally treat the case of deep networks, which we have omitted here for simplicity. The explicit formula for the NTK of ReLU networks as presented in Examples 11.17 and 11.31 was given in [36]. The observation that neural networks at initialization behave like Gaussian processes presented in Section 11.5.5 was first made in [155]. For a general reference on Gaussian processes see the textbook [186]. When only training the last layer of a network (in which the network is affine linear), there are strong links to random feature methods [184]. Recent developments on this topic can also be found in the literature under the name “Neural network Gaussian processes”, or NNGPs for short [128, 46].

Exercises

Exercise 11.32. Prove Theorem 11.3.

Hint: Assume first that $\mathbf{w}_0 \in \ker(\mathbf{A})^\perp$ (i.e. $\mathbf{w}_0 \in \tilde{H}$). For $\text{rank}(\mathbf{A}) < d$, using $\mathbf{w}_k = \mathbf{w}_{k-1} - h\nabla f(\mathbf{w}_{k-1})$ and the singular value decomposition of \mathbf{A} , write down an explicit formula for \mathbf{w}_k . Observe that due to $1/(1-x) = \sum_{k \in \mathbb{N}_0} x^k$ for all $x \in (0, 1)$ it holds $\mathbf{w}_k \rightarrow \mathbf{A}^\dagger \mathbf{y}$ as $k \rightarrow \infty$, where \mathbf{A}^\dagger is the Moore-Penrose pseudoinverse of \mathbf{A} .

Exercise 11.33. Let $\mathbf{x}_i \in \mathbb{R}^d$, $i = 1, \dots, m$. Show that there exists a “feature map” $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$, such that for any configuration of labels $y_i \in \{-1, 1\}$, there always exists a hyperplane in \mathbb{R}^m separating the two sets $\{\phi(\mathbf{x}_i) \mid y_i = 1\}$ and $\{\phi(\mathbf{x}_i) \mid y_i = -1\}$.

Exercise 11.34. Consider the RBF kernel $K : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, $K(x, x') := \exp(-(x - x')^2)$. Find a Hilbert space H and a feature map $\phi : \mathbb{R} \rightarrow H$ such that $K(x, x') = \langle \phi(x), \phi(x') \rangle_H$.

Exercise 11.35. Let $n \in \mathbb{N}$ and consider the polynomial kernel $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^n$. Find a Hilbert space H and a feature map $\phi : \mathbb{R}^d \rightarrow H$, such that $K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_H$.

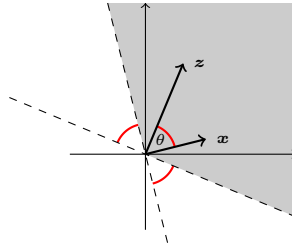
Hint: Use the multinomial formula.

Exercise 11.36. Let $u_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$ be i.i.d. standard Gaussian distributed random variables for $i = 1, \dots, d$. Show that for all nonzero $\mathbf{x}, \mathbf{z} \in \mathbb{R}^d$

$$\mathbb{E}[\mathbf{1}_{[0, \infty)}(\mathbf{u}^\top \mathbf{x}) \mathbf{1}_{[0, \infty)}(\mathbf{u}^\top \mathbf{z})] = \frac{\pi - \theta}{2\pi}, \quad \theta = \arccos\left(\frac{\mathbf{x} \mathbf{z}^\top}{\|\mathbf{x}\| \|\mathbf{z}\|}\right).$$

This shows the formula for the ReLU NTK with Gaussian initialization as discussed in Example 11.31.

Hint: Consider the following sketch



Exercise 11.37. Consider the network (11.5.1) with LeCun initialization as in (11.5.3), but with the biases instead initialized as

$$c, b_i \stackrel{\text{iid}}{\sim} \mathcal{D}(0, 1) \quad \text{for all } i = 1, \dots, n. \quad (11.6.4)$$

Compute the corresponding NTK as in Theorem 11.23. Moreover, compute the NTK also for the normalized network (11.6.1) with initialization (11.6.2) as in Theorem 11.30, but replace again the bias initialization with that given in (11.6.4).