

Deep learning foundations: a unified analysis of optimization, generalization and representation

Ziwei Ji

ziweiji2@illinois.edu

University of Illinois, Urbana-Champaign

Contents

- 1 Introduction** **2**

- 2 Preliminaries** **4**
 - 2.1 Supervised learning 4
 - 2.2 Function classes 5
 - 2.3 Empirical risk minimization 5
 - 2.4 Gradient flow, gradient descent, and stochastic gradient descent 6

- 3 The neural tangent kernel** **6**
 - 3.1 Regression 7
 - 3.2 Binary classification 9
 - 3.2.1 Representation 9
 - 3.2.2 Optimization 10
 - 3.2.3 Generalization 11
 - 3.3 Open problems 12

- 4 The implicit bias and margin maximization** **12**
 - 4.1 Linear classifiers, separable case 13
 - 4.2 Linear classifiers, nonseparable case 15
 - 4.3 Deep homogeneous networks 15
 - 4.4 Open problems 17

- 5 Other open problems** **17**

1 Introduction

Deep neural networks have achieved great empirical successes in recent years. They form the backbone of many state-of-the-art algorithms in computer vision, natural language processing and reinforcement learning, which further have many applications in object detection, machine translation, self-driving cars, etc.

Meanwhile, theoretical analysis of neural networks still lags far behind practice. Important questions are not fully answered in all core parts of learning theory:

- Optimization: our ultimate goal is to find a model which has a small test error, i.e., a good accuracy on future/unseen data. However, we cannot minimize the test error directly; typically we have a training set, and first minimize the training error. In deep learning, the training error is usually highly nonconvex, but standard optimization algorithms such as gradient descent usually work pretty well, even though they have no convergence guarantee in a general nonconvex setting. *Why is training feasible in deep learning?*
- Generalization: suppose we can achieve a low training error, the next step is to ensure a low generalization error, i.e., a small difference between the training error and test error. There are many classical tools to bound the generalization error, such as the VC dimension bound: it depends on the number of parameters in a network, and will be small if the network can only represent a limited number of sign patterns. However, in practice, the number of weights in a network is usually much larger than the number of training examples, and the network can even fit random labels (Zhang et al., 2016). Therefore we cannot directly apply the VC dimension bound to explain the good generalization of deep networks. *How to give a fine-grained analysis of the good generalization of deep networks?*
- Representation: it is well-established that a neural network with only one hidden layer and a non-polynomial activation can already approximate any continuous function (Leshno et al., 1993). In practice, deep networks can also usually achieve zero training error. However, classical representation results usually suffer from the curse of dimension: to achieve an approximation error of ϵ with input dimension d , in the worse case $\Omega(1/\epsilon^d)$ nodes are needed. In addition, even if there exists a good representation, it is still unclear if it can be found by a practical algorithm. *How to avoid the curse of dimension, and how to build a representation theory within the range reachable by the training algorithm?*

To attack the above problems, we suggest studying optimization, generalization, and representation of deep networks simultaneously. In this manuscript, we will summarize the progress made so far, and also discuss a few potential future directions. Most discussion will be focused on gradient descent (GD) and stochastic gradient descent (SGD), since they are already very successful in practice despite their simplicity, and lots of research has been done on them. On the other hand, there are many other interesting and efficient algorithms, which will also be briefly discussed.

Here we briefly describe the methodologies and results. Roughly speaking, the training process of a deep network using gradient descent have an early phase and a late phase.

The early phase. This phase is often handled using the *neural tangent kernel* (Jacot et al., 2018), or the NTK, which is very influential in recent years. In this analysis, the first observation is that the *linearized model at initialization* (i.e., the first-order Taylor approximation around the initialization) is already strong enough to give zero training error (Du et al., 2018c), as long as the network is wide enough. More precisely, if we use a *linear model* with the *random features* given by the initialization, we can already fit any training data. The optimization and generalization analysis for this linearized model is also relatively simple, since linear models are well-studied. Specifically, the test error bound would depend on the *inverse NTK margin* (cf. Section 3.2.3).

The next observation is that, once again with a wide enough network, if we also choose a suitable initialization and learning rate, then the gradient descent path on the true network will be close to the gradient descent path on the linearized model, for a period that is long enough to guarantee a low training error and even low test error (Li and Liang, 2018; Du et al., 2018c,b; Allen-Zhu et al., 2018b; Zou et al., 2018; Allen-Zhu et al., 2018a; Arora et al., 2019b). The pitfall is that, these analyses usually require the width of the network to be polynomial in the number of training examples, or the inverse failure probability, or the inverse target error; such a large width is never used in practice.

In (Ji and Telgarsky, 2019b), we prove that for binary classification, it is actually enough to have the width depend *polylogarithmically* on all of the aforementioned parameters; the width only needs to depend polynomially on the inverse NTK margin, which is necessary if we want to use an NTK-style analysis. We also provide sample complexity bounds for GD and SGD, and moreover prove that our SGD sample complexity bound is *tight* in the NTK regime, which suggests that our analysis captures the full power of the NTK in some cases.

The late phase. The NTK analysis requires the weights to stay close to initialization. For a fixed number of training iterations, this can be true if the network is wide enough. However, in practice people keep training as long as their computational budget allows (Shallue et al., 2018); in this setting, the norm of weights will keep increasing to infinity with the cross entropy loss, and therefore the underlying assumption of the NTK cannot hold forever. Moreover, as shown in (Arora et al., 2019a), finite-width deep networks usually work better than the corresponding infinite-width networks/NTK. Consequently, finite-width neural networks must have some special properties beyond the NTK.

Motivated by the empirical success of gradient descent, we study the *implicit bias* of gradient descent: that is, among all low-training-error solutions, we try to prove that gradient descent prefers “simple” solutions, and thus generalize well on unseen data. In particular, we prove that in some cases, gradient descent can find the *global maximum margin*, which could be much larger than the NTK margin. A test error bound then follows from standard margin-based bounds (Bartlett et al., 2017).

Here is a summary of recent results on the implicit bias and margin maximization of gradient descent:

- Soudry et al. (2018) prove that with linear classifiers, linearly separable data and exponentially-tailed losses, gradient descent converges in direction to the ℓ_2 -maximum-margin direction.

- Still with linear classifiers, in (Ji and Telgarsky, 2018b), we further characterize the implicit bias of gradient descent for general nonseparable data, and in (Ji and Telgarsky, 2019a), we show that a dual-based analysis can further characterize the implicit bias for general non-exponentially-tailed losses, and establish faster rates for margin maximization and convergence to the implicit bias.
- Surprisingly, margin maximization still holds for deep linear networks (Gunasekar et al., 2018b; Ji and Telgarsky, 2018a) and 2-homogeneous networks (Chizat and Bach, 2020), and can be generalized in a certain form to deep homogeneous networks which allows ReLU, max-pooling, etc. (Lyu and Li, 2019; Ji and Telgarsky, 2020), even though the objective function is nonconvex in such settings. In particular, we prove *directional convergence* and *alignment* for general homogeneous models in (Ji and Telgarsky, 2020), which generalize most prior implicit bias results.

Here is an outline of this manuscript. In Section 2, we introduce the problem setting. In Section 3, we focus on the neural tangent kernel, and discuss existing results and open directions. The implicit bias and margin maximization properties of gradient descent are covered in Section 4. In Section 5, we explore other interesting open directions.

2 Preliminaries

2.1 Supervised learning

In this manuscript, we focus on supervised learning. Other models such as unsupervised learning and reinforcement learning are also very interesting to study, but we will not discuss them in detail here.

We are given a feature space $\mathcal{X} \subset \mathbb{R}^d$, a label space \mathcal{Y} (with concrete examples given below), and a joint distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$. We want to learn a function $f : \mathcal{X} \rightarrow \widehat{\mathcal{Y}}$ which matches the distribution \mathcal{D} , where $\widehat{\mathcal{Y}}$ denotes the output space of f , which may or may not be the same as \mathcal{Y} (see examples below). Formally, we have a loss function $\ell : \widehat{\mathcal{Y}} \times \mathcal{Y} \rightarrow \mathbb{R}$, and we try to minimize the *risk* of f , defined as

$$\mathcal{R}(f) := \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f(x), y)].$$

Here are some concrete examples:

- Regression: here we have $\mathcal{Y} = \widehat{\mathcal{Y}} = \mathbb{R}^m$, and a commonly-used loss is the *squared loss*:

$$\ell(f(x), y) := \frac{1}{2} \|f(x) - y\|_2^2.$$

- Classification: here the label space is a finite space $\mathcal{Y} = \{1, 2, \dots, m\}$, while $\widehat{\mathcal{Y}}$ is usually \mathbb{R}^m , with the *cross entropy loss*:

$$\ell(f(x), y) := -\ln \left(\frac{\exp(f(x)_y)}{\sum_{j=1}^m \exp(f(x)_j)} \right).$$

Other loss functions and learning problems are discussed in Section 5.

2.2 Function classes

In this manuscript, we try to find a risk-minimizing f which lies in a certain parameterized function class \mathcal{F} . For example, the class of linear classifiers is parameterized by a weight vector $w \in \mathbb{R}^d$:

$$\mathcal{F}_{\text{lin}} := \left\{ x \mapsto \langle w, x \rangle \mid w \in \mathbb{R}^d \right\},$$

while the class of deep fully-connected neural networks is parameterized by weight matrices W_k and bias vectors b_k :

$$\mathcal{F}_{\text{nn}} := \left\{ x \mapsto W_L \sigma_{L-1}(\cdots W_2 \sigma_1(W_1 x + b_1) + b_2 \cdots) + b_L \right\},$$

where $W_k \in \mathbb{R}^{m_k \times d_k}$ (particularly, $d_1 = d$), and $b_k \in \mathbb{R}^{m_k}$, and $\sigma_k : \mathbb{R}^{m_k} \rightarrow \mathbb{R}^{d_{k+1}}$ are (nonlinear) activation functions, such as ReLU, max-pooling, etc. Note that \mathcal{F}_{lin} and \mathcal{F}_{nn} defined here are quite general, and we will almost always add more constraints. Other models, such as convolutional networks and residual networks, will also be discussed.

2.3 Empirical risk minimization

To find a good f , we sample a training set $S = \{(x_i, y_i)\}_{i=1}^n$ from the distribution \mathcal{D} , and try to minimize the *empirical risk*:

$$\widehat{\mathcal{R}}_S(f) := \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i).$$

For simplicity, we often just use $\widehat{\mathcal{R}}(f)$. With the notion of empirical risk, we can give a formal description of the problems of representation, optimization and generalization: given a training set S , suppose our learning algorithm finds a solution f_S . Note that

$$\mathcal{R}(f_S) = \left(\mathcal{R}(f_S) - \widehat{\mathcal{R}}_S(f_S) \right) + \left(\widehat{\mathcal{R}}_S(f_S) - \widehat{\mathcal{R}}_S(f^*) \right) + \left(\widehat{\mathcal{R}}_S(f^*) - \mathcal{R}(f^*) \right) + \mathcal{R}(f^*),$$

where f^* is an arbitrary fixed member of \mathcal{F} . If we can ensure all four parts on the right hand side are small, then $\mathcal{R}(f_S)$ is small.

1. The first part gives the problem of generalization. It looks similar to the third part; however, as $|S| \rightarrow \infty$, the third part converges to 0 due to the law of large numbers, while it is unclear whether the first part also converges to 0 since f_S depends on S .
2. The second part corresponds to optimization; it is small if we can minimize the empirical risk on \mathcal{F} .
3. The fourth part is called representation; it basically asks whether the function class is powerful enough to minimize the true risk.

Representation, optimization and generalization are often studied independently. However, as discussed in the introduction, to build a satisfactory theory for deep learning, it looks necessary to analyze them together. For example, an optimization analysis usually requires a representation result, and gives some norm and margin bounds which can be used to prove a generalization bound.

2.4 Gradient flow, gradient descent, and stochastic gradient descent

To minimize the empirical risk, it is usually enough to use first-order methods. Suppose the function/predictor f is parameterized by a weight vector w , i.e., f is a function of the input x and weight vector w , and we can write $f_w(x) = f(x, w)$, and

$$\widehat{\mathcal{R}}(w) = \widehat{\mathcal{R}}(f_w) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, w), y_i).$$

The simplest algorithm is gradient descent (GD): we start from some initialization w_0 , and for any $t \geq 0$, let

$$w_{t+1} := w_t - \eta_t \nabla \widehat{\mathcal{R}}(w_t), \tag{1}$$

where $\eta_t > 0$ is the step size / learning rate. When $\eta_t \rightarrow 0$, i.e., when we consider infinitesimal step sizes, then we get gradient flow: it is the solution to the following differential equation

$$\frac{dw_t}{dt} = -\nabla \widehat{\mathcal{R}}(w_t). \tag{2}$$

Gradient flow can usually simplify the analysis a lot and illustrate the main ideas. However, we should be cautious when applying gradient-flow results to gradient descent: the second-order error term in the discrete-time case can sometimes be hard to control.

In the literature, stochastic gradient descent (SGD) may have two different meanings: the first one assumes that at step t , we have a fresh sample $(x_t, y_t) \sim \mathcal{D}$, and do

$$w_{t+1} := w_t - \eta_t \left. \frac{\partial \ell(f(x_t, w), y_t)}{\partial w} \right|_{w=w_t}.$$

The second meaning of SGD, which is also what people use in practice, assumes that we have a training set S , and in each epoch of training, S is randomly decomposed into batches S_1, \dots, S_B , and we just do B steps of GD where $\widehat{\mathcal{R}}_{S_b}$ is used in the b -th step. In this manuscript we will use the first meaning of SGD, but it is very interesting and important to handle the second SGD.

Other training algorithms, such as the momentum and adaptive gradient methods, will be discussed in Section 5.

3 The neural tangent kernel

Let us first discuss the intuition behind the neural tangent kernel. For simplicity, suppose the network has a 1-dimensional output.

Note that with linear predictors and convex loss functions, the empirical risk $\widehat{\mathcal{R}}$ is also convex, and there are many tools from convex optimization to analyze it. However, deep networks are highly nonlinear functions, and $\widehat{\mathcal{R}}$ is highly nonconvex in this setting, which makes it hard to analyze. To solve this issue, one way is to consider the *linearization* of

a network around its initialization: given a classifier $f(w; x)$ where w denotes the weight vector and x denotes the input, let w_0 denote the initialization, then by Taylor’s theorem,

$$f(w; x) \approx f(w_0; x) + \langle \nabla f(w_0; x), w - w_0 \rangle. \quad (3)$$

The right hand side of eq. (3) is basically the *neural tangent model*. It also has a kernel formulation, which will be introduced in Section 3.1, called the *neural tangent kernel*; for simplicity, in the following we will mostly use “NTK” to refer to this linearization idea.

An NTK analysis of GD roughly consists of the following steps: (i) show that in the overparameterized regime (i.e., when the width is large), with high probability over the Gaussian initialization of w_0 , on the linear model with features given by $\nabla f(w_0; x)$, GD can minimize the empirical risk to 0; (ii) prove that with overparameterization, GD iterates on the neural network are close to GD iterates on the linear model in (i).

Below we discuss the above ideas in more detail on a two-layer network. As in (Du et al., 2018c; Arora et al., 2019b; Ji and Telgarsky, 2019b), the network has weight matrices $W \in \mathbb{R}^{m \times d}$ and $a \in \mathbb{R}^m$, and outputs

$$f(x; W, a) := \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \sigma(\langle w_s, x \rangle), \quad (4)$$

with initialization

$$w_{s,0} \sim \mathcal{N}(0, I_d), \quad \text{and} \quad a_s \sim \text{unif}(\{-1, +1\}),$$

where $w_{s,t}$ denotes the s -th row of W at step t . (It turns out that this specific initialization and the $1/\sqrt{m}$ factor induce the NTK behavior (Chizat et al., 2019).) We consider the ReLU activation $\sigma(z) := \max\{0, z\}$, though the analyses can be easily extended to Lipschitz continuous, positively homogeneous activations such as the leaky ReLU. (We can define the derivative of ReLU at 0 to be an arbitrary number in $[0, 1]$, and results in this section will hold.) For simplicity, we also make the technical assumptions that a is fixed and only W is trained, and $\|x_i\|_2 = 1$ for all feature vector x_i in the training set, as in (Li and Liang, 2018; Du et al., 2018c; Arora et al., 2019b; Ji and Telgarsky, 2019b).

3.1 Regression

Du et al. (2018c) consider a regression problem and the squared loss. Formally, given a training set $\{(x_i, y_i)\}_{i=1}^n$, and a network as in eq. (4), we try to minimize the empirical risk

$$\widehat{\mathcal{R}}(W) := \frac{1}{2n} \sum_{i=1}^n (f(x_i; W, a) - y_i)^2.$$

The key tool used in the regression analysis is the neural tangent kernel: the infinite-width NTK on the training set is an n -by- n matrix, where the (i, j) -th entry is defined as

$$H^\infty(i, j) := \langle x_i, x_j \rangle \mathbb{E}_{w \sim \mathcal{N}(0, I_d)} \left[\mathbf{1} [\langle w, x_i \rangle > 0, \langle w, x_j \rangle > 0] \right].$$

Although this expression is complicated, it is basically the inner product between gradients at initialization, as can be seen from the corresponding finite-width NTK:

$$\begin{aligned} H_{W_0}(i, j) &:= \left\langle \left. \frac{\partial f(x_i; W, a)}{\partial W} \right|_{W_0}, \left. \frac{\partial f(x_j; W, a)}{\partial W} \right|_{W_0} \right\rangle \\ &= \frac{1}{m} \langle x_i, x_j \rangle \sum_{s=1}^m \mathbb{1} [\langle w_{s,0}, x_i \rangle > 0, \langle w_{s,0}, x_j \rangle > 0]. \end{aligned}$$

Note that H_{W_0} is a random matrix with expectation H^∞ . Du et al. (2018c) assume that the smallest eigenvalue of H^∞ is positive.

Assumption 3.1. It holds that $\lambda_0 := \lambda_{\min}(H^\infty) > 0$.

Assumption 3.1 can be viewed as a *representation* assumption on the training data: it means that we can achieve zero training error using the infinite-width NTK features at initialization. Assumption 3.1 is true as long as there are no parallel feature vectors in the training set (Du et al., 2018c, Theorem 3.1).

With Assumption 3.1, Du et al. (2018c) prove the following result.

Theorem 3.1. (Du et al., 2018c, Theorem 4.1) *Under Assumption 3.1, with the number of hidden units $m = \Omega\left(\frac{n^6}{\lambda_0^4 \delta^3}\right)$ and step size $\eta = O\left(\frac{\lambda_0}{n}\right)$, with probability $1 - \delta$ over the random initialization,*

$$\widehat{\mathcal{R}}(W_t) \leq \left(1 - \frac{\eta \lambda_0}{2n}\right)^t \widehat{\mathcal{R}}(W_0).$$

Here we briefly discuss the proof ideas, which is clearer for gradient flow. For the squared loss, we can prove the following results, which demonstrates why the NTK is natural in this setting:

$$\frac{d\widehat{\mathcal{R}}(W_t)}{dt} \leq -\frac{2\lambda_{\min}(H_{W_t})}{n} \widehat{\mathcal{R}}(W_t). \quad (5)$$

If we can further show that, for example, it always holds that $\lambda_{\min}(H_{W_t}) \geq \lambda_0/2$, then the convergence rate follows from Grönwall's inequality.

Du et al. (2018c) then consider a suitably chosen ρ and the following set given W_0 :

$$\mathcal{W}_\rho := \left\{ W \in \mathbb{R}^{m \times d} \mid \|w_s - w_{s,0}\|_2 \leq \rho \text{ for all } 1 \leq s \leq m \right\}. \quad (6)$$

The proof idea is to show that for all $W \in \mathcal{W}_\rho$, it holds that $\lambda_{\min}(H_W) \geq \lambda_0/2$. Therefore before gradient flow exits \mathcal{W}_ρ , the empirical risk $\widehat{\mathcal{R}}$ decreases very fast, as ensured by eq. (5). The final step is to use eq. (5) and Grönwall's inequality to show that gradient flow actually never exits \mathcal{W}_ρ . This is possible as long as the width is large enough, since the chosen ρ is independent of the width m , but the movement of each w_s in gradient flow is roughly $1/\sqrt{m}$.

Is polynomial overparameterization necessary? The previous analysis relies upon the squared loss, and the required width depends polynomially on n , the number of training examples, and $1/\delta$, the inverse failure probability. Arora et al. (2019b) extend the above analysis and give a generalization bound, which also require the squared loss and at least the same amount of overparameterization. Oymak and Soltanolkotabi (2019); Song and Yang (2019) further reduce the required overparameterization, but there is still a $\text{poly}(n)$ dependency. In fact, since these results do not assume any relationship between features and labels, a $\text{poly}(n)$ dependency might be inevitable.

On the other hand, if there is a strong feature-label relationship, then it is reasonable to expect a much milder dependency on n . As discussed in the next subsection, in (Ji and Telgarsky, 2019b), we prove that it is indeed possible to get a *polylogarithmic* dependency on n and other parameters, using the logistic loss and the notion of *NTK separation margin*.

3.2 Binary classification

Here we discuss our results in (Ji and Telgarsky, 2019b). We consider binary classification, meaning that the label y_i is either $+1$ or -1 . We use the logistic loss $\ell(z) := \ln(1 + e^{-z})$, and

$$\widehat{\mathcal{R}}(W) := \frac{1}{n} \sum_{i=1}^n \ell(y_i f(x_i; W, a)).$$

One property we have not fully exploited so far is the *homogeneity* of the network: let $f_i(W) := f(x_i; W, a)$. Then for any W , it holds that $\langle \nabla f_i(W), W \rangle = f_i(W)$. (This is one version of Euler’s homogeneous function theorem.) In other words, the network is basically a linear classifier with the gradients as features, but these features may change during training.

3.2.1 Representation

Similar to regression, here we also need to make a representation assumption at initialization. Instead of positive eigenvalues, we want a positive *separation margin* at initialization, i.e., with high probability, there exists $\bar{U} \in \mathbb{R}^{m \times d}$ which can separate $\{(\nabla f_i(W_0), y_i)\}_{i=1}^n$ with a positive margin:

$$\min_{1 \leq i \leq n} \left(y_i \langle \bar{U}, \nabla f_i(W_0) \rangle \right) = \min_{1 \leq i \leq n} \left(y_i \frac{1}{\sqrt{m}} \sum_{s=1}^m a_s \langle \bar{u}_s, x_i \rangle \mathbb{1} [\langle w_{s,0}, x_i \rangle > 0] \right) > 0. \quad (7)$$

To get a deterministic version of eq. (7), below we state its infinite-width limit, with an additional norm bound on the separator. Let $\mu_{\mathcal{N}}$ denote the Gaussian measure on \mathbb{R}^d , given by the Gaussian density with respect to the Lebesgue measure on \mathbb{R}^d . We consider the following Hilbert space

$$\mathcal{H} := \left\{ w : \mathbb{R}^d \rightarrow \mathbb{R}^d \mid \int \|w(z)\|_2^2 d\mu_{\mathcal{N}}(z) < \infty \right\}.$$

For any $x \in \mathbb{R}^d$, define $\phi_x \in \mathcal{H}$ by

$$\phi_x(z) := x \mathbb{1} [\langle z, x \rangle > 0],$$

and particularly define $\phi_i := \phi_{x_i}$ for the training input x_i .

Assumption 3.2. There exists $\bar{v} \in \mathcal{H}$ and $\gamma > 0$, such that $\|\bar{v}(z)\|_2 \leq 1$ for any $z \in \mathbb{R}^d$, and for any $1 \leq i \leq n$,

$$y_i \langle \bar{v}, \phi_i \rangle_{\mathcal{H}} := y_i \int \langle \bar{v}(z), \phi_i(z) \rangle d\mu_{\mathcal{N}}(z) \geq \gamma.$$

The space \mathcal{H} is the reproducing kernel Hilbert space (RKHS) induced by the infinite-width NTK, and ϕ_x maps x into \mathcal{H} . Assumption 3.2 supposes that the induced training set $\{(\phi_i, y_i)\}_{i=1}^n$ can be separated by some $\bar{v} \in \mathcal{H}$, with an additional bound on $\|\bar{v}(z)\|_2$ which is crucial in our analysis: specifically, \bar{v} will be used to construct a separator \bar{U} as in eq. (7), upon which our optimization analysis is centered.

This assumption, and the distributional version Assumption 3.3, are introduced in (Nittanda and Suzuki, 2019) for smooth activations. (Cao and Gu, 2019) make a similar separability assumption, but in the RKHS induced by the second layer a ; by contrast, Assumption 3.3 is on separability in the RKHS induced by the first layer W .

How large can γ be? An interesting and important question is how large γ could be. As shown below in Theorems 3.2 to 3.4, although our required width depends polylogarithmically on every other parameter, it still depends polynomially on $1/\gamma$. If γ is very small, we may still need a strong overparameterization.

In fact, γ can be small in general: in (Ji and Telgarsky, 2019b, Proposition 5.2), we prove that it is $1/\text{poly}(n)$ with random labels. However, it actually makes sense, since we should not expect a large margin for random labels.

On the other hand, consider the noisy 2-XOR example (Wei et al., 2018), where each coordinate of the feature vector is an i.i.d. Rademacher random variable, and the label is the XOR of the first two coordinates. We prove that even though the training set can have 2^d examples, the margin $\gamma = \Omega(1/d) = \Omega(1/\ln(n))$ (Ji and Telgarsky, 2019b, Proposition 5.3). In fact, this margin is large enough to help us prove a tight sample complexity bound (see Section 3.2.3). This suggests that the above notion of margin can capture the full power of the NTK around initialization in some cases.

3.2.2 Optimization

Under Assumption 3.2, we prove the following result.

Theorem 3.2. (Ji and Telgarsky, 2019b, Theorem 2.2) *Under Assumption 3.2, given any risk target $\epsilon \in (0, 1)$ and any $\delta \in (0, 1/3)$, let*

$$\lambda := \frac{\sqrt{2 \ln(4n/\delta)} + \ln(4/\epsilon)}{\gamma/4}, \quad \text{and} \quad M := \frac{4096\lambda^2}{\gamma^6}.$$

Then for any $m \geq M$ and any constant step size $\eta \leq 1$, with probability $1 - 3\delta$ over the random initialization,

$$\frac{1}{T} \sum_{t < T} \widehat{\mathcal{R}}(W_t) \leq \epsilon, \quad \text{where} \quad T := \lceil 2\lambda^2/\eta\epsilon \rceil.$$

Moreover for any $0 \leq t < T$ and any $1 \leq s \leq m$,

$$\|w_{s,t} - w_{s,0}\|_2 \leq \frac{4\lambda}{\gamma\sqrt{m}}.$$

As shown above, our required width is $\text{polylog}(n, 1/\delta, 1/\epsilon) \cdot \text{poly}(1/\gamma)$. The most closely related work is (Nitanda and Suzuki, 2019), which considers smooth activations and requires $\tilde{\Omega}(1/\epsilon^2)$ hidden units. Li and Liang (2018) consider SGD with the cross entropy loss and two-layer networks, and need $\text{poly}(l, 1/\epsilon)$ hidden units, where l is a quantity depends on the data distribution. Allen-Zhu et al. (2018a) consider SGD on a two-layer network, and a variant of SGD on a three-layer network; they assume a ground truth network with infinite-order smooth activations, and their required width depends polynomially on $1/\epsilon$ and some constants related to the smoothness of the activations of the ground truth network. Cao and Gu (2019) consider deep networks, and need $\Omega(1/\epsilon^{14})$ hidden units. Following our work, Chen et al. (2019) prove a polylogarithmic width for deep networks.

Is a $\text{poly}(1/\gamma)$ dependency necessary? It is natural to ask whether it is possible to also get a $\text{polylog}(1/\gamma)$ dependency. However, this is impossible if we want to use an NTK analysis: in (Ji and Telgarsky, 2019b, Proposition 5.7), we construct a training set with 4 examples, such that if $m \leq \sqrt{d} - 2/4$ for $d \geq 20$, then with a constant probability over the random initialization of W_0 , it holds that $\{(\nabla f_i(W_0), y_i)\}_{i=1}^4$ is nonseparable. Consequently, we cannot use an NTK analysis to prove something like Theorem 3.2 in this setting.

3.2.3 Generalization

To get a generalization bound, we extends Assumption 3.2 to the data distribution \mathcal{D} .

Assumption 3.3. There exists $\bar{v} \in \mathcal{H}$ and $\gamma > 0$, such that $\|\bar{v}(z)\|_2 \leq 1$ for any $z \in \mathbb{R}^d$, and

$$y \int \langle \bar{v}(z), x \rangle \mathbf{1}[\langle z, x \rangle > 0] d\mu_{\mathcal{N}}(z) \geq \gamma$$

for almost all (x, y) sampled from the data distribution \mathcal{D} .

Here is the test error bound we prove.

Theorem 3.3. (Ji and Telgarsky, 2019b, Theorem 3.2) Under Assumption 3.3, given any $\epsilon, \delta \in (0, 1)$, using a constant step size no larger than 1 and let

$$n = \tilde{\Omega}\left(\frac{1}{\gamma^4\epsilon^2}\right), \quad \text{and} \quad m = \Omega\left(\frac{\ln(n/\delta) + \ln(1/\epsilon)^2}{\gamma^8}\right),$$

it holds with probability $1 - \delta$ that $P_{(x,y) \sim \mathcal{D}}(yf(x; W_k, a) \leq 0) \leq \epsilon$, where k denotes the step with the minimum empirical risk in the first $\tilde{\Theta}(1/\gamma^2\epsilon)$ steps.

The proof uses a Rademacher complexity bound based on the sigmoid function $z \mapsto 1/(1 + e^z)$, which is also equal to $-\ell'(z)$ for the logistic loss. Similar to regression, the proof also uses some \mathcal{W}_ρ (cf. eq. (6)), where ρ is chosen based on the bound on $\|w_{s,t} - w_{s,0}\|_2$ given by Theorem 3.2.

Using similar proof ideas as above, and a martingale Bernstein bound, we further prove the following test error bound for SGD.

Theorem 3.4. *Under Assumption 3.3, given any $\epsilon, \delta \in (0, 1)$, using a constant step size and $m = \Omega\left(\frac{(\ln(1/\delta) + \ln(1/\epsilon^2))}{\gamma^8}\right)$, it holds with probability $1 - \delta$ that*

$$\frac{1}{n} \sum_{i=1}^n P_{(x,y) \sim \mathcal{D}}(yf(x; W_i, a) \leq 0) \leq \epsilon, \quad \text{for } n = \tilde{\Theta}(1/\gamma^2 \epsilon).$$

A tight sample complexity bound. Although in (Ji and Telgarsky, 2019b), we focus on two-layer network, all our analyses work for the infinite-width NTK. In particular, we can prove an identical sample complexity bound for SGD as Theorem 3.4. For the noisy 2-XOR example, since $\gamma = \Omega(1/d)$, it follows that to achieve a constant test accuracy we need $\tilde{O}(d^2)$ samples. On the other hand, (Wei et al., 2018) give a d^2 sample complexity lower bound for any NTK classifier on the noisy 2-XOR data. Therefore our SGD sample complexity upper bound is *tight* up to a logarithmic factor.

3.3 Open problems

There are many interesting open problems to investigate:

- Deep networks: The idea of NTK has been used a lot to analyze deep networks (Du et al., 2018b; Allen-Zhu et al., 2018b; Zou et al., 2018). However, a benefit of depth has not been shown in the NTK framework. Can we prove a better separation margin for deep networks?
- Convolutional networks: the same question as above can be asked for convolutional networks. We also expect a better margin here in light of the recent work (Arora et al., 2019a), which empirically shows that the convolutional neural tangent kernel achieves a test accuracy on CIFAR which is close to the performance of the corresponding finite deep networks.
- Beyond the initial phase: the linearization idea eq. (3) can be applied around any weights, such as W_t . However, the NTK analysis works by requiring $\nabla f_i(W_t) \approx \nabla f_i(W_0)$, in which case we cannot show that training a neural network is a better idea than training the corresponding NTK. Can we show that the NTK, or the gradient features $\nabla f_i(W_t)$, are improving during training?

4 The implicit bias and margin maximization

In the previous section, using an NTK analysis, we show that GD can minimize the empirical risk to an arbitrarily low level, and we also prove a test error bound which depends on the

inverse NTK margin. However, in the separable case, with the logistic loss, GD iterates will go to the infinity, and thus will eventually exit the early NTK phase. It means that we cannot use the NTK analysis after a long training, even though practitioners run their optimization methods as long as their computational budget allows (Shallue et al., 2018). Moreover, the NTK margin is usually not the global maximum margin, and can be small.

In this section, we focus on the late phase of training, and try to analyze the *implicit bias* of gradient descent: among all the low-training-error solutions, which one GD prefers. In particular, we will focus on binary classification, and prove *global margin maximization* in a few cases; generalization bounds which depend on the inverse global maximum margin then follow from standard margin-based bounds (Bartlett et al., 2017).

4.1 Linear classifiers, separable case

Here we consider linear classifiers, and assume that training set is linearly separable.

Assumption 4.1. There exists $u \in \mathbb{R}^d$ such that $y_i \langle u, x_i \rangle > 0$ for all i .

For simplicity, we will focus on the exponential loss:

$$\ell(z) := e^{-z}, \quad \text{and} \quad \widehat{\mathcal{R}}(w) := \frac{1}{n} \sum_{i=1}^n \ell(y_i \langle w, x_i \rangle) = \frac{1}{n} \sum_{i=1}^n \exp(-y_i \langle w, x_i \rangle).$$

However, results presented here can also be extended to other exponentially-tailed losses, such as the logistic loss (with more complicated proofs though).

The implicit bias: ℓ_2 margin maximization. Under Assumption 4.1, it holds that $\inf_{w \in \mathbb{R}^d} \widehat{\mathcal{R}}(w) = 0$; however this infimum is never attained, since the exponential function never attains its infimum 0. Consequently, the optimal w is *off at infinity*, and if we run GD on this problem, we have $\lim_{t \rightarrow \infty} \|w_t\|_2 = \infty$.

Although w_t does not converge, it is still possible for its *direction* $w_t / \|w_t\|_2$ to converge. In fact, for binary classification, a good direction is all we need. Indeed, Soudry et al. (2018) prove that for exponentially-tailed losses, the GD iterates converge in direction to the ℓ_2 -maximum-margin direction \bar{u} :

$$\bar{\gamma} := \max_{\|w\|_2 \leq 1} \min_{1 \leq i \leq n} y_i \langle w, x_i \rangle, \quad \text{and} \quad \bar{u} := \arg \max_{\|w\|_2 \leq 1} \min_{1 \leq i \leq n} y_i \langle w, x_i \rangle.$$

Naturally, the ℓ_2 margin is also maximized:

$$\lim_{t \rightarrow \infty} \frac{\min_{1 \leq i \leq n} y_i \langle w_t, x_i \rangle}{\|w_t\|_2} = \bar{\gamma}.$$

A margin maximization proof based on primal smoothness. Since the analysis of directional convergence (i.e., the convergence of $w_t / \|w_t\|_2$) is more complicated, here we give a proof sketch for margin maximization (Telgarsky, 2013; Gunasekar et al., 2018a; Ji and Telgarsky, 2018b). This proof can also show directional convergence, but can only give a suboptimal rate.

The proof is centered upon a *smoothed (unnormalized) margin* (Lyu and Li, 2019), which is also called a *generalized sum* (Hardy et al., 1934):

$$\psi(w) := \ell^{-1} \left(\sum_{i=1}^n \ell(y_i \langle w, x_i \rangle) \right).$$

For the exponential loss, we have

$$\psi(w) = \ell^{-1} \left(n \widehat{\mathcal{R}}(w_t) \right) \leq \min_{1 \leq i \leq n} y_i \langle w_t, x_i \rangle \leq \psi(w) + \ln(n), \quad (8)$$

and therefore $\psi(w)$ is indeed an approximation of the true unnormalized margin. It is also 1-smooth with respect to the ℓ_∞ norm, and thus called the smoothed unnormalized margin.

To prove margin maximization, the key step is to show the following:

$$\frac{\psi(w_{t+1}) - \psi(w_t)}{\|w_{t+1}\|_2 - \|w_t\|_2} \geq \bar{\gamma} - o(1). \quad (9)$$

Margin maximization then follows from eqs. (8) and (9) and that $\|w_t\|_2 \rightarrow \infty$. The denominator of eq. (9) is handled using the triangle inequality, while the numerator can be handled using the ℓ_∞ smoothness of ψ ; this analysis can give a $\ln(t)/\sqrt{t}$ margin maximization rate with a suitably chosen step size schedule, and can actually handle general steepest descent such as boosting. Next we show that a faster $O(1/t)$ rate, specifically for gradient descent, can be obtained using a dual analysis.

Faster convergence rates via dual multiplicative weight updates. In (Ji and Telgarsky, 2019a), we show that (primal) gradient descent on the empirical risk induces a multiplicative weight update (or mirror descent, or dual averaging) on a dual objective. From this dual perspective, the “implicit” margin maximization property of GD is actually explicit: the optimum of the dual objective is exactly given by the maximum margin. More interestingly, using the smoothness of the dual objective, we can prove a faster margin maximization and implicit bias rate.

The dual variable $q_t \in \Delta_n$ is given by the softmax mapping on the unnormalized margin distribution:

$$q_{t,i} := \frac{\exp(-y_i \langle w_t, x_i \rangle)}{\sum_{j=1}^n \exp(-y_j \langle w_t, x_j \rangle)}.$$

Moreover, the dual objective is defined as

$$g(q) := \frac{1}{2} \left\| \sum_{i=1}^n y_i q_i x_i \right\|_2^2. \quad (10)$$

It can be verified that the update from q_t to q_{t+1} is exactly a mirror descent update on the dual objective g , and moreover the dual iterates q_t can minimize g over Δ_n (Ji and Telgarsky, 2019a, Theorem 2.2). Note that the minimum of g is $\bar{\gamma}^2/2$, and moreover the optimal solution

\bar{q} satisfies $\sum_{i=1}^n y_i \bar{q}_i x_i = \bar{\gamma} \bar{u}$ (Ji and Telgarsky, 2018b, Theorem 2.1). Consequently the dual iterates *explicitly* maximize the margin.

Additionally, the dual objective g is 1-smooth with respect to the ℓ_1 norm; this allows us to prove a faster convergence rate for $\psi(w_t)$ than merely using the smoothness of ψ . Once again using eq. (9), we can prove the following faster margin maximization rate.

Theorem 4.1. (Ji and Telgarsky, 2019a, Theorem 4.2) *Under Assumption 4.1, for the exponential loss, if $w_0 = 0$ and $\eta_t = 1/\widehat{\mathcal{R}}(w_t)$, then*

$$\frac{\min_{1 \leq i \leq n} y_i \langle w_t, x_i \rangle}{\|w_t\|} \geq \bar{\gamma} - \frac{\ln(n) + 1}{\bar{\gamma} t}.$$

4.2 Linear classifiers, nonseparable case

In (Ji and Telgarsky, 2018b), we still consider linear classifiers, but make no assumption on the training set, i.e., the training data may or may not be linearly separable. We prove that we can decompose the training set S into a separable and a nonseparable part in a unique way, on which GD has different implicit biases.

Formally, given a dataset $S = \{(x_i, y_i)\}_{i=1}^n$, we decompose it into $S_{\text{sep}} \cup S_{\text{sc}}$ in the following way: for each data example (x_i, y_i) , if there exists a unit vector u such that $y_i \langle u, x_i \rangle > 0$ and $y_j \langle u, x_j \rangle \geq 0$ for all $1 \leq j \leq n$, then we include (x_i, y_i) into S_{sep} , otherwise we include it into S_{sc} . Define

$$\widehat{\mathcal{R}}_{\text{sep}}(w) := \frac{1}{n} \sum_{(x_i, y_i) \in S_{\text{sep}}} \ell(y_i \langle w, x_i \rangle), \quad \text{and} \quad \widehat{\mathcal{R}}_{\text{sc}}(w) := \frac{1}{n} \sum_{(x_i, y_i) \in S_{\text{sc}}} \ell(y_i \langle w, x_i \rangle),$$

and $H := \text{span}(\{x_i : (x_i, y_i) \in S_{\text{sc}}\})$. It follows that $\widehat{\mathcal{R}} = \widehat{\mathcal{R}}_{\text{sep}} + \widehat{\mathcal{R}}_{\text{sc}}$, and that S_{sep} is linearly separable by some linear classifier on H^\perp , and that $\widehat{\mathcal{R}}_{\text{sc}}$ is strongly convex on H , and thus has a unique minimizer \bar{v} on H .

In (Ji and Telgarsky, 2018b), we further prove that, without knowing this decomposition at all, GD on $\widehat{\mathcal{R}}$ can find the unique minimizer \bar{v} on H , and converges in direction to the maximum margin separator of S_{sep} on H^\perp .

4.3 Deep homogeneous networks

Linear classifier is the most simple case of a *positively homogeneous* model, which is used heavily throughout the theoretical study of deep networks (Du et al., 2018a; Lyu and Li, 2019; Woodworth et al., 2020). Let $f(x; w)$ denote a model with input x and parameter w , we say that it is L -positively homogeneous if for all $c > 0$, it holds that $f(x; cw) = c^L f(x; w)$. In this subsection, we will discuss our implicit bias results in (Ji and Telgarsky, 2020), which hold for all well-behaved positively-homogeneous model and can imply most prior margin maximization results.

For a network to be “well-behaved”, in addition to positive homogeneity, we also need the network to be *definable in an o-minimal structure*. This is a mild condition, and allows most common layers used in deep learning. More specifically, our analysis can handle ReLU

and leaky ReLU activations, max-pooling and average-pooling, fully-connected layers, convolutional layers, etc. We cannot handle skip connections and bias vectors, since they are excluded by homogeneity; it is an interesting open problem to handle them.

We also assume that at initialization, we already have $\widehat{\mathcal{R}}(w_0) < 1/n$. This is also assumed in the prior work (Lyu and Li, 2019), and can be ensured using an NTK analysis. By making this assumption, we explicitly focus on the late phase of training.

We prove the following results for gradient flow (Ji and Telgarsky, 2020, Theorem 3.1 and 4.1):

- The gradient flow iterate w_t converges in direction:

$$\lim_{t \rightarrow \infty} \frac{w_t}{\|w_t\|_2} \text{ exists.}$$

- Gradient flow cannot converge to an arbitrary direction; it must also be satisfied that $-\nabla \widehat{\mathcal{R}}(w_t)$ converges to the same direction.

$$\lim_{t \rightarrow \infty} \left\langle \frac{w_t}{\|w_t\|_2}, \frac{-\nabla \widehat{\mathcal{R}}(w_t)}{\|\nabla \widehat{\mathcal{R}}(w_t)\|_2} \right\rangle = 1.$$

Many implicit bias works explicitly assume directional convergence and some version of alignment (Gunasekar et al., 2018b; Chizat and Bach, 2020), but neither do these works indicate a possible proof, nor do they provide conclusive evidence.

Directional convergence is proved using unbounded nonsmooth Kurdyka-Łojasiewicz inequalities. Prior works on similar topics cannot be applied here since they either assume a real analytic objective function (Kurdyka et al., 2000) and therefore cannot handle ReLU or max-pooling, or excludes the exponential function (Kurdyka et al., 2006; Grandjean, 2007) and thus cannot handle the exponential or logistic loss.

Alignment is proved by extending the dual objective g (cf. eq. (10)) to the homogeneous setting, and showing that it still converges.

Margin maximization consequences. As a sanity check, we can show that the implicit bias for linear classifiers already follows from alignment. Here are other margin maximization consequences for deep networks:

- For deep linear networks, margin maximization immediately follows from directional convergence and alignment (Ji and Telgarsky, 2020, Proposition 4.4). Previously, this is proved by Gunasekar et al. (2018b) assuming directional convergence of weights and gradients; Ji and Telgarsky (2018a) removed such assumptions, but need to assume the support vectors span \mathbb{R}^d . Our proof is cleaner and requires fewer assumptions.
- For two-homogeneous networks, global margin maximization is proved by Chizat and Bach (2020) for infinite-width networks, assuming directional convergence and one version of dual convergence. Using directional convergence and alignment, we give a cleaner proof of the finite-width case, using a covering condition taken from the infinite-width analysis (Ji and Telgarsky, 2020, Proposition 4.6).

4.4 Open problems

One interesting open problem is to see if GD also maximizes the margin for deep homogeneous networks and convolutional networks. Gunasekar et al. (2018b) characterizes the implicit bias of deep linear convolutional networks; however their result only shows a first-order stationary point of the margin, and cannot handle nonlinear activations. It is interesting to see if margin maximization still follows from directional convergence and alignment in a more general setting.

Another interesting open problem is to prove a good test error bound, even *without* exact margin maximization. For example, for linear classifiers, the $1/t$ margin maximization rate given in Theorem 4.1 can only be achieved using an aggressive step size schedule; if we use a constant step size with the empirical risk (i.e., a constant η_t in eq. (1)), then the margin maximization rate is $\Theta(\ln(n)/\ln(t))$, which is very slow. However, in this setting, Shamir (2020) still proves an $\tilde{O}\left(\frac{1}{\bar{\gamma}^2 t} + \frac{1}{\bar{\gamma}^2 n}\right)$ test error bound, which does not require exact margin maximization. It is interesting to see if we can prove a similar test error bound for deep networks.

5 Other open problems

There are many other interesting open problems; here is an incomplete list of them.

Other architectures. Although fully-connected layers are still widely used, other components such as convolutional layers, skip connections, batch normalization layers, etc. are also very popular in modern deep networks; it is therefore very interesting to theoretically analyze them.

For example, Li et al. (2020) shows that fully-connected networks, trained by common algorithms including gradient descent, satisfy a property called orthogonal equivariance. They also construct an example on which there exists a convolutional network that is more sample efficient on any orthogonal equivariant model. It is interesting to see if orthogonal equivariance can explain the good generalization of convolutional networks in practice, and whether there are other explanations.

Other loss functions. Typically, the squared loss is used for regression problems, while the logistic loss and cross entropy loss are used for classification problems. However, Hui and Belkin (2020) show that the squared loss also works well in practice for classification tasks. How to explain this, and can we design better loss functions?

Other training algorithms. In practice, SGD is usually run with *momentum*. Moreover, algorithms with *adaptive learning rates*, such as Adam (Kingma and Ba, 2014), are also very popular. However, it is believed that adaptive gradient methods learn models with poor generalization (Keskar and Socher, 2017; Wilson et al., 2017). When does this happen, and can we improve the generalization of adaptive gradient methods?

References

- Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. *arXiv preprint arXiv:1811.04918*, 2018a.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. *arXiv preprint arXiv:1811.03962*, 2018b.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. *arXiv preprint arXiv:1904.11955*, 2019a.
- Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. *arXiv preprint arXiv:1901.08584*, 2019b.
- Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural networks. In *Advances in Neural Information Processing Systems*, pages 6240–6249, 2017.
- Yuan Cao and Quanquan Gu. Generalization error bounds of gradient descent for learning over-parameterized deep relu networks. *arXiv preprint arXiv:1902.01384*, 2019.
- Zixiang Chen, Yuan Cao, Difan Zou, and Quanquan Gu. How much over-parameterization is sufficient to learn deep relu networks? *arXiv preprint arXiv:1911.12360*, 2019.
- Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. *arXiv preprint arXiv:2002.04486*, 2020.
- Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. In *Advances in Neural Information Processing Systems*, pages 2937–2947, 2019.
- Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. In *Advances in Neural Information Processing Systems*, pages 384–395, 2018a.
- Simon S Du, Jason D Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global minima of deep neural networks. *arXiv preprint arXiv:1811.03804*, 2018b.
- Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018c.
- V Grandjean. On the limit set at infinity of a gradient trajectory of a semialgebraic function. *Journal of Differential Equations*, 233(1):22–41, 2007.
- Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. *arXiv preprint arXiv:1802.08246*, 2018a.

- Suriya Gunasekar, Jason D Lee, Daniel Soudry, and Nati Srebro. Implicit bias of gradient descent on linear convolutional networks. In *Advances in Neural Information Processing Systems*, pages 9461–9471, 2018b.
- G. H. Hardy, J. E. Littlewood, and G. Pólya. *Inequalities*. Cambridge university press, 1934.
- Like Hui and Mikhail Belkin. Evaluation of neural architectures trained with square loss vs cross-entropy in classification tasks. *arXiv preprint arXiv:2006.07322*, 2020.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580, 2018.
- Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. *arXiv preprint arXiv:1810.02032*, 2018a.
- Ziwei Ji and Matus Telgarsky. Risk and parameter convergence of logistic regression. *arXiv preprint arXiv:1803.07300*, 2018b.
- Ziwei Ji and Matus Telgarsky. Characterizing the implicit bias via a primal-dual analysis. *arXiv preprint arXiv:1906.04540*, 2019a.
- Ziwei Ji and Matus Telgarsky. Polylogarithmic width suffices for gradient descent to achieve arbitrarily small test error with shallow relu networks. *arXiv preprint arXiv:1909.12292*, 2019b.
- Ziwei Ji and Matus Telgarsky. Directional convergence and alignment in deep learning. *arXiv preprint arXiv:2006.06657*, 2020.
- Nitish Shirish Keskar and Richard Socher. Improving generalization performance by switching from adam to sgd. *arXiv preprint arXiv:1712.07628*, 2017.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Krzysztof Kurdyka, Tadeusz Mostowski, and Adam Parusinski. Proof of the gradient conjecture of r. thom. *Annals of Mathematics*, 152(3):763–792, 2000.
- Krzysztof Kurdyka, Adam Parusiński, et al. Quasi-convex decomposition in o-minimal structures. application to the gradient conjecture. In *Singularity theory and its applications*, pages 137–177. Mathematical Society of Japan, 2006.
- Moshe Leshno, Vladimir Ya Lin, Allan Pinkus, and Shimon Schocken. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural networks*, 6(6):861–867, 1993.
- Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in Neural Information Processing Systems*, pages 8157–8166, 2018.

- Zhiyuan Li, Yi Zhang, and Sanjeev Arora. Why are convolutional nets more sample-efficient than fully-connected nets? *arXiv preprint arXiv:2010.08515*, 2020.
- Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. *arXiv preprint arXiv:1906.05890*, 2019.
- Atsushi Nitanda and Taiji Suzuki. Refined generalization analysis of gradient descent for over-parameterized two-layer neural networks with smooth activations on classification problems. *arXiv preprint arXiv:1905.09870*, 2019.
- Samet Oymak and Mahdi Soltanolkotabi. Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. *arXiv preprint arXiv:1902.04674*, 2019.
- Christopher J. Shallue, Jaehoon Lee, Joseph Antognini, Jascha Sohl-Dickstein, Roy Frostig, and George E. Dahl. Measuring the effects of data parallelism on neural network training. 2018. [arXiv:1811.03600 \[cs.LG\]](#).
- Ohad Shamir. Gradient methods never overfit on separable data. *arXiv preprint arXiv:2007.00028*, 2020.
- Zhao Song and Xin Yang. Quadratic suffices for over-parametrization via matrix chernoff bound. *arXiv preprint arXiv:1906.03593*, 2019.
- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.
- Matus Telgarsky. Margins, shrinkage, and boosting. In *ICML*, 2013.
- Colin Wei, Jason D Lee, Qiang Liu, and Tengyu Ma. Regularization matters: Generalization and optimization of neural nets vs their induced kernel. *arXiv preprint arXiv:1810.05369*, 2018.
- Ashia C Wilson, Rebecca Roelofs, Mitchell Stern, Nati Srebro, and Benjamin Recht. The marginal value of adaptive gradient methods in machine learning. In *Advances in neural information processing systems*, pages 4148–4158, 2017.
- Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. *arXiv preprint arXiv:2002.09277*, 2020.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes over-parameterized deep relu networks. *arXiv preprint arXiv:1811.08888*, 2018.