

Training neural networks: Dynamics and convergence

Instructor: Lei Wu

PKU Summer School, 2021

Outline

- Classical convergence results of gradient descent (GD).
- Hardness in training neural networks.
- Convergence in the over-parameterized regime.

A brief review of classical results for gradient descent

Gradient descent

- Consider the problem of minimizing

$$\min_{\theta} \hat{\mathcal{R}}_n(\theta).$$

- The GD iterates as follows

$$\theta_{t+1} = \theta_t - \eta_t \nabla \hat{\mathcal{R}}_n(\theta_t),$$

where η_t is the learning rate.

- When $\eta_t \rightarrow 0$, the GD becomes the GD flow:

$$\frac{d\theta_t}{dt} = -\nabla \hat{\mathcal{R}}_n(\theta_t).$$

Convergence of GD

Theorem 1 (Non-convex)

For any $t > 0$,

$$\min_{s \in [0, t]} \|\nabla \hat{\mathcal{R}}_n(\theta_s)\| \leq \sqrt{\frac{\hat{\mathcal{R}}_n(\theta_0) - \inf_{\theta} \hat{\mathcal{R}}_n(\theta)}{t}}.$$

Theorem 2

Assume that $\hat{\mathcal{R}}_n$ is convex and the minimizer is given by θ^* with $\|\theta^*\|_2 < \infty$. Then, we have

$$\hat{\mathcal{R}}_n(\theta_t) - \hat{\mathcal{R}}_n(\theta^*) \leq \frac{\|\theta^* - \theta_0\|_2^2}{2t}.$$

Convergence of GD (Cont'd)

- Can we prove the converge to global minima for non-convex problem? This problem often strongly depends on the specific model. There exists a general condition as follows.
- $\hat{\mathcal{R}}_n$ is said to satisfy the Polyak-Lojasiewicz (PL) condition if

$$\|\nabla \hat{\mathcal{R}}_n(\theta)\|_2 \geq C(\hat{\mathcal{R}}_n(\theta) - \inf_{\theta} \hat{\mathcal{R}}_n(\theta)).$$

Theorem 3

Under the PL condition, we have

$$\hat{\mathcal{R}}_n(\theta_t) - \inf_{\theta} \hat{\mathcal{R}}_n(\theta) \leq e^{-Ct}(\hat{\mathcal{R}}_n(\theta_0) - \inf_{\theta} \hat{\mathcal{R}}_n(\theta)).$$

**Hardness in learning Barron functions:
Training two-layer neural networks suffers
CoD**

Motivation

- Both the approximation and estimation error for Barron functions obey the Monte-Carlo rate, which is free of CoD. A nature questions is then: *Do there exist algorithms that can learn these functions efficiently?*

Motivation

- Both the approximation and estimation error for Barron functions obey the Monte-Carlo rate, which is free of CoD. A natural question is then: *Do there exist algorithms that can learn these functions efficiently?*
- We say an algorithm \mathcal{A} is efficient in learning a function class \mathcal{F} , if for every $\varepsilon > 0$, $f^* \in \mathcal{F}$, the time complexity of returning a solution \hat{f} such that $\|\hat{f} - f^*\| \leq \varepsilon$ satisfies:

Time complexity = $\text{poly}(1/\varepsilon, d)$.

Motivation

- Both the approximation and estimation error for Barron functions obey the Monte-Carlo rate, which is free of CoD. A natural question is then: *Do there exist algorithms that can learn these functions efficiently?*
- We say an algorithm \mathcal{A} is efficient in learning a function class \mathcal{F} , if for every $\varepsilon > 0$, $f^* \in \mathcal{F}$, the time complexity of returning a solution \hat{f} such that $\|\hat{f} - f^*\| \leq \varepsilon$ satisfies:

$$\text{Time complexity} = \text{poly}(1/\varepsilon, d).$$

- We will show that the class of Barron functions is not efficiently learnable.

Learning intersections of halfspaces

- Let $x \in \mathcal{X} = \{-1, 1\}^d$ and consider the binary classification problem, i.e., $f^* : \mathcal{X} \mapsto \{-1, 1\}$.
- Let σ_{step} be the step function, i.e., $\sigma_{\text{step}}(t) = 1(t \geq 0)$.

Learning intersections of halfspaces

- Let $x \in \mathcal{X} = \{-1, 1\}^d$ and consider the binary classification problem, i.e., $f^* : \mathcal{X} \mapsto \{-1, 1\}$.
- Let σ_{step} be the step function, i.e., $\sigma_{\text{step}}(t) = 1(t \geq 0)$.

We will need the following hardness result for learning the intersection of halfspaces.

Theorem 4 (Theorem 1.2, Kalai, Klivans, 2008)

Let $\mathcal{H} = \{x \mapsto \sigma_{\text{step}}(w^T x - b - 1/2) : b \in \mathbb{N}, w \in \mathbb{N}^d, |b| + \|w\|_1 \leq \text{poly}(d)\}$. Define

$$\mathcal{H}_K = \{x \mapsto h_1(x) \wedge h_2(x) \wedge \cdots \wedge h_K(x) : h_i \in \mathcal{H}\}.$$

Assume $k \geq d^\rho$ with $\rho > 0$. Then, under a certain **cryptographic** assumption, \mathcal{H}_K is not efficiently learnable.

- The proof is to reduce it to some classical hard problems, e.g., k -coloring. The **cryptographic** assumption means that we assume that these hard problems are indeed hard in certain sense. If this assumption does not hold, the modern cryptosystem can be broken in a polynomial time.
- We shall show two-layer neural networks can simulate the functions in \mathcal{H}_K .

Hardness of learning two-layer ReLU networks

Theorem 5 (Livni, et al, 2014)

Let $\mathcal{X} = \{-1, 1\}^d$, and $\mathcal{G} = \{f \in \mathcal{B} : \|f\|_{\mathcal{B}} \leq \text{poly}(d)\}$. Then, \mathcal{G} is not efficiently learnable.

- The intuition is that 2-layer neural network can simulate the intersections of hyperspaces.
- The step function can be approximated by two ReLU functions very well:

$$\sigma_{\text{step}}(t) = \lim_{a \rightarrow \infty} \text{ReLU}(at) - \text{ReLU}(at - 1).$$

Hardness of learning two-layer neural networks

Proof:

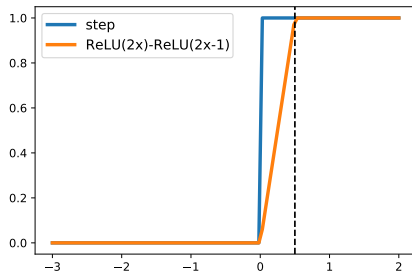
- Let $c(x) = w^T x - b - 1/2$. Since $w \in \mathbb{N}^d, x \in \{-1, 1\}^d, b \in \mathbb{N}$, we have $|c(x)| \geq 1/2$. Assume $\|w\|_1 + |b| \leq \text{poly}(d)$.
- Consider k hyperplanes $\{c_i\}_{i=1}^k$. Let $h_i(x) = \sigma_{\text{step}}(c_i(x)) \in \mathcal{H}$.
- Let

$$g(x) = \frac{1}{2k} \left(\sum_{i=1}^k (\text{ReLU}(2c_i(x)) - \text{ReLU}(2c_i(x) - 1)) - k + \frac{1}{3} \right).$$

Obviously, g is a 2-layer ReLU network with the path norm bounded by

$$\frac{1}{2k} \left(k + \frac{1}{3} + \sum_{i=1}^k (2(2\|w_i\|_1 + |b_i| + 1/2) + 1) \right) = \text{poly}(d).$$

- The blue part is equal to $\sigma_{\text{step}}(c_i(x))$ due to $\sigma_{\text{step}}(z) = \text{ReLU}(2z) - \text{ReLU}(2z - 1)$ for $|z| \geq 1/2$.



- We can verify that

$$\text{sign}(g(x)) = h_1(x) \wedge h_2(x) \wedge \cdots \wedge h_k(x), \quad \forall x \in \{-1, 1\}^d.$$

- Note that similar results also hold for two-layer networks with the sigmoid activation function, since the sigmoid function can approximate the step function as well. See [Livni, et al, 2014] for more details.
- The above results rely on the hardness of certain classical hard problems.
 - Pros: It is implied that the hardness holds for any algorithms.
 - Cons: This perspective is too abstract. It does not provide any concrete examples and intuitions behind the hardness of training.
- In the following, we will provide some understandings from a landscape perspective.

Orthonormal classes

Denote by \mathcal{D} the distribution over the input space \mathcal{X} . For any two functions f_1, f_2 , define the inner product $\langle f_1, f_2 \rangle = \mathbb{E}_{x \sim \mathcal{D}}[f_1(x)f_2(x)]$.

Definition 6 (Orthonormal class)

Let \mathcal{F} be a function class. We say that it is an orthonormal class, if $\langle f_i, f_j \rangle = \delta_{i,j}$ for any $f_i, f_j \in \mathcal{F}$.

Orthonormal classes

Denote by \mathcal{D} the distribution over the input space \mathcal{X} . For any two functions f_1, f_2 , define the inner product $\langle f_1, f_2 \rangle = \mathbb{E}_{x \sim \mathcal{D}}[f_1(x)f_2(x)]$.

Definition 6 (Orthonormal class)

Let \mathcal{F} be a function class. We say that it is an orthonormal class, if $\langle f_i, f_j \rangle = \delta_{i,j}$ for any $f_i, f_j \in \mathcal{F}$.

- Let $\mathcal{B}_d = \{f \in \mathcal{B} : \|f\|_{\mathcal{B}} \lesssim d^2\}$. We will show that \mathcal{B}_d contains an orthonormal subset $\mathcal{F} = \{f_1, \dots, f_m\}$ with $m = \exp(d)$.
- We will show that learning the orthonormal class \mathcal{F} is hard if $|\mathcal{F}| = \exp(d)$.

Parity functions

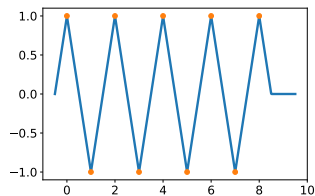
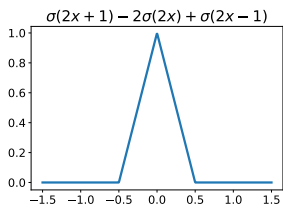
- $\mathcal{F}_1 = \{f_v(x) = (-1)^{\langle v, x \rangle} : v \in \{0, 1\}^d\}$, where $x \in \{0, 1\}^d$.
- Consider $\mathcal{D} = \text{Unif}(\{0, 1\}^d)$. Then, we have

$$\langle f_v, f_{v'} \rangle = \mathbb{E}_x [(-1)^{(v+v')^T x}] = \mathbb{E}_x \left[\prod_{i=1}^d (-1)^{(v_i+v'_i)x_i} \right] \quad (0.1)$$

$$= \prod_{i=1}^d \mathbb{E}_{x_i} [(-1)^{(v_i+v'_i)x_i}] = \delta_{v, v'}. \quad (0.2)$$

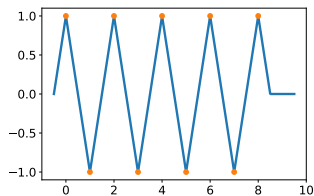
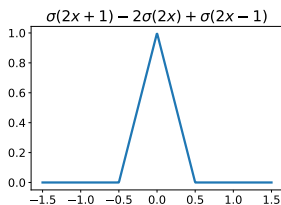
Hence, \mathcal{F}_1 is an orthonormal class with $|\mathcal{F}_1| = 2^d$.

Parity functions as two-layer neural networks



- Observation:
 - The function $(-1)^s$ for $s \in \mathbb{N}$ can be implemented using the triangle wave.
 - The triangle wave can be written as a linear combination of the hat function, which is a linear combination of ReLU function (left figure).

Parity functions as two-layer neural networks



- Observation:
 - The function $(-1)^s$ for $s \in \mathbb{N}$ can be implemented using the triangle wave.
 - The triangle wave can be written as a linear combination of the hat function, which is a linear combination of ReLU function (left figure).
- Note that $v^T x \in \{0, 1, \dots, d\}$. Let σ be the ReLU function. Then,

$$(-1)^{v^T x} = \sigma_{tri}(v^T x) = \sum_{i=0}^d (-1)^i \sigma_{hat}(v^T x - i) \quad (0.3)$$

$$= \sum_{i=0}^d (-1)^i (\sigma(2v^T x - 2i + 1) - 2\sigma(2v^T x - 2i) + \sigma(2v^T x - 2i - 1)) \quad (0.4)$$

- It is easy to show that the path norm of this network is bounded by Cd^2 .

Cosine neurons

Consider the domain $\mathcal{X} = [0, 2\pi]^d$.

- Let $S_d = \{\cos(w^T x) : w \in \mathbb{N}^d, \sum_{i=1}^d w_i \leq d\}$.
- In the previous lecture, we have shown that

$$|S_d| \gtrsim \frac{2^d}{\sqrt{d}} \quad (0.5)$$

$$\|f\|_{\mathcal{B}} \lesssim d^2, \quad \forall f \in S_d. \quad (0.6)$$

- Hence, for this continuous case, \mathcal{B}_d contains exponential many orthonormal functions.

Remark

- $\mathcal{F}_\phi := \{\phi(w^T x) : \|w\| = \sqrt{d}\}$. $\mathcal{D} = \mathcal{N}(0, I_d)$. [Shamir, 2017] shows that as long as ϕ is periodic, under some mild condition, \mathcal{F}_ϕ contains an orthonormal subset $\mathcal{F}_1 = \{f_1, \dots, f_m\}$ with $m = \exp(d)$.

Gradients for an orthonormal class

Let $h(\cdot; \theta)$ be any parametric model. Denote by $R^f(\theta) = \mathbb{E}_x[(h(x; \theta) - f(x))^2]$ the risk. Then, we have the following theorem.

Theorem 7

Let \mathcal{F} be an orthonormal class. Let P denote the uniform distribution over the space of \mathcal{F} and $g(\theta) = \mathbb{E}_{f \sim P}[\nabla R^f(\theta)]$. We have

$$\mathbb{E}_{f \sim P}[(\nabla R^f(\theta) - g(\theta))^2] \leq \frac{\mathbb{E}_x[\|\nabla_{\theta} h(x; \theta)\|^2]}{|\mathcal{F}|}. \quad (0.7)$$

Gradients for an orthonormal class

Let $h(\cdot; \theta)$ be any parametric model. Denote by $R^f(\theta) = \mathbb{E}_x[(h(x; \theta) - f(x))^2]$ the risk. Then, we have the following theorem.

Theorem 7

Let \mathcal{F} be an orthonormal class. Let P denote the uniform distribution over the space of \mathcal{F} and $g(\theta) = \mathbb{E}_{f \sim P}[\nabla R^f(\theta)]$. We have

$$\mathbb{E}_{f \sim P}[(\nabla R^f(\theta) - g(\theta))^2] \leq \frac{\mathbb{E}_x[\|\nabla_{\theta} h(x; \theta)\|^2]}{|\mathcal{F}|}. \quad (0.7)$$

- If $|\mathcal{F}|$ is exponentially in d , e.g., the parity functions. The variance of gradients is exponentially small.

Gradients for an orthonormal class

Let $h(\cdot; \theta)$ be any parametric model. Denote by $R^f(\theta) = \mathbb{E}_x[(h(x; \theta) - f(x))^2]$ the risk. Then, we have the following theorem.

Theorem 7

Let \mathcal{F} be an orthonormal class. Let P denote the uniform distribution over the space of \mathcal{F} and $g(\theta) = \mathbb{E}_{f \sim P}[\nabla R^f(\theta)]$. We have

$$\mathbb{E}_{f \sim P}[(\nabla R^f(\theta) - g(\theta))^2] \leq \frac{\mathbb{E}_x[\|\nabla_{\theta} h(x; \theta)\|^2]}{|\mathcal{F}|}. \quad (0.7)$$

- If $|\mathcal{F}|$ is exponentially in d , e.g., the parity functions. The variance of gradients is exponentially small.
- This theorem implies that the “information” about the target function contained in the gradient is exponentially small. Therefore, one would expect that gradient-based methods will be unlikely to learn the function class \mathcal{F} .

Proof:

- First, the gradient can be written as follows

$$\nabla_{\theta} R^f = \mathbb{E}_x[(h(x; \theta) - f) \nabla_{\theta} h(x; \theta)] = C_{\theta} - \langle f, \nabla_{\theta} h(x; \theta) \rangle,$$

where C_{θ} is independent of the target function f .

Proof:

- First, the gradient can be written as follows

$$\nabla_{\theta} R^f = \mathbb{E}_x[(h(x; \theta) - f) \nabla_{\theta} h(x; \theta)] = C_{\theta} - \langle f, \nabla_{\theta} h(x; \theta) \rangle,$$

where C_{θ} is independent of the target function f .

- Hence,

$$\mathbb{E}_f[(\nabla_{\theta} R^f - g(\theta))^2] \leq \mathbb{E}_f[\langle f, \nabla_{\theta} h(x; \theta) \rangle^2] \tag{0.8}$$

$$\leq \frac{1}{|\mathcal{F}|} \sum_f \langle f, \nabla_{\theta} h(x; \theta) \rangle^2 \tag{0.9}$$

$$\leq \frac{\mathbb{E}_x[\|\nabla_{\theta} h(x; \theta)\|^2]}{|\mathcal{F}|}. \tag{0.10}$$

Hardness of learning with GD: Setup

Setup:

- Assume \mathcal{F} to be an orthonormal class with $|\mathcal{F}| = 2^d$. Consider the binary classification with the hinge loss. The risk is given by

$$R^f(\theta) := \mathbb{E}_x[\max(0, 1 - h(x; \theta)f(x))]. \quad (0.11)$$

- Assume $|h(x; \theta)| \leq 1$ and $|f(x)| \leq 1$ for any $x \in \mathcal{X}$. Then we have

$$\begin{aligned} \mathbb{E}_f[\|\nabla_\theta R^f(\theta)\|^2] &= \mathbb{E}_f(\mathbb{E}_x[f(x)\nabla_\theta h(x; \theta)])^2 \\ &= \frac{1}{|\mathcal{F}|} \sum_i \langle f_i, \nabla h(\cdot; \theta) \rangle^2 \leq \frac{\|\nabla_\theta h\|^2}{|\mathcal{F}|} \leq \frac{G_\theta}{2^d}. \end{aligned}$$

Remark: the above assumption holds for parity functions.

Hardness of learning with GD

Theorem 8

Assume the model satisfies that $\sup_{x \in X} |h(x; \theta)| \leq 1$ and $\mathbb{E}_x[\|\nabla_{\theta} h(x; \theta_1) - \nabla_{\theta} h(x; \theta_2)\|^2] \leq L\|\theta_1 - \theta_2\|^2$. Let θ_0, θ_t^f be the GD solution at time 0 and time t , respectively. Then, there exist C_1, C_2 such that

$$\mathbb{E}_f[\|\theta_t^f - \theta_0\|^2] \leq C_1(e^{\frac{C_2 t}{2^{d/2}}} - 1), \quad (0.12)$$

where C_1, C_2 only depend on L and θ_0 .

Hardness of learning with GD

Theorem 8

Assume the model satisfies that $\sup_{x \in X} |h(x; \theta)| \leq 1$ and $\mathbb{E}_x [\|\nabla_{\theta} h(x; \theta_1) - \nabla_{\theta} h(x; \theta_2)\|^2] \leq L \|\theta_1 - \theta_2\|^2$. Let θ_0, θ_t^f be the GD solution at time 0 and time t , respectively. Then, there exist C_1, C_2 such that

$$\mathbb{E}_f [\|\theta_t^f - \theta_0\|^2] \leq C_1 (e^{\frac{C_2 t}{2^{d/2}}} - 1), \quad (0.12)$$

where C_1, C_2 only depend on L and θ_0 .

The above theorem implies that GD solution is exponentially close to the initialization in polynomial time. More rigorously, we have the following corollary.

Corollary 9

For any $T = \text{poly}(d)$, there exists a $f \in \mathcal{F}$ such that

$$\|\theta_t^f - \theta_0\| \leq C \frac{\text{poly}(d)}{2^d}, \quad \forall t \in [0, T]$$

where C only depends on L and θ_0 .

Hardness of learning with GD

Proof:

- $G(\theta) = \mathbb{E}_x[\|\nabla_{\theta} h(x; \theta)\|^2]$ satisfies

$$G(\theta) \leq G(\theta_0) + 2L\|\theta - \theta_0\|^2. \quad (0.13)$$

- Therefore,

$$\frac{d \mathbb{E}_f[\|\theta_t^f - \theta_0\|^2]}{dt} = 2 \mathbb{E}_f[\langle \theta_t^f - \theta_0, -\nabla_{\theta} R^f(\theta_t^f) \rangle] \quad (0.14)$$

$$\leq \frac{1}{2^{\frac{d}{2}-1}} \sqrt{\mathbb{E}_f[\|\theta_t^f - \theta_0\|^2] \mathbb{E}_f[G(\theta_t^f)]} \quad (0.15)$$

$$\leq \frac{1}{2^{\frac{d}{2}-1}} \sqrt{\mathbb{E}_f[\|\theta_t^f - \theta_0\|^2] \mathbb{E}_f[G(\theta_0) + 2L\|\theta_t^f - \theta_0\|^2]}. \quad (0.16)$$

Hardness of learning with GD

Proof: Let $\delta_t = \sqrt{\mathbb{E}_f[\|\theta_t^f - \theta_0\|^2]}$. Then, we have

$$\dot{\delta}_t \leq 2^{2-\frac{d}{2}}(\sqrt{2L}\delta_t + \sqrt{G(\theta_0)}). \quad (0.17)$$

By Gronwall's inequality, we obtain

$$\delta_t \leq \sqrt{\frac{G(\theta_0)}{2L}}(e^{2^{2-\frac{d}{2}}\sqrt{L}t} - 1).$$

Numerical evidence

Consider learning parity functions with online SGD. Fig. 1 shows the convergence of SGD. Here, the model is two-layer neural networks with width being 2000. The hinge loss $\ell(y, y') = \max(0, 1 - yy')$ is used, batch size is 2000 and learning rate is 0.002. We see clearly that when $d = 20$, the training process does not show any improvement in a reasonable time.

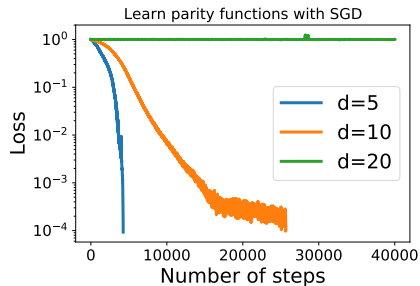
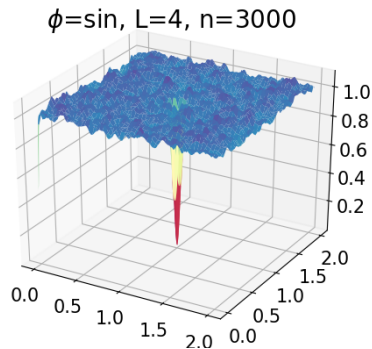
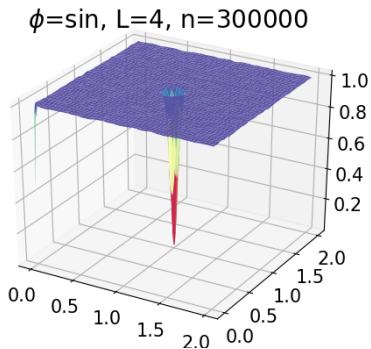


Figure 1: Learning parity functions with SGD and two-layer neural networks.

An illustration of the landscape of cosine neuron

Consider a two-dimensional case. $R_L(w) = \mathbb{E}_{x \sim \mathcal{D}}[(\sin(Lw^T x) - \sin(Lw^{*T} x))^2]$, where L can be viewed as a proxy of the dimension d .



- For the population landscape, the global minima locate in a deep well with other place is extremely flat. This confirms Theorem 5.
- The empirical landscape is full of bad local minima.

Summary

- Learning a subset of two-layer neural networks, whose path norms are bounded by $\text{poly}(d)$, can be reduced to certain classical hard problems, whose hardness is assumed to be true. Otherwise, the modern cryptosystem can be broken in polynomial time. This type of hardness results hold for any algorithms.

Summary

- Learning a subset of two-layer neural networks, whose path norms are bounded by $\text{poly}(d)$, can be reduced to certain classical hard problems, whose hardness is assumed to be true. Otherwise, the modern cryptosystem can be broken in polynomial time. This type of hardness results hold for any algorithms.
- For orthonormal classes, we show that the gradient variance (wrt the target function) of is exponentially small. Hence, *gradient-based* algorithms are unlikely to succeed. This observation hold for any parametric model as long as they satisfy some smooth condition.

Summary

- Learning a subset of two-layer neural networks, whose path norms are bounded by $\text{poly}(d)$, can be reduced to certain classical hard problems, whose hardness is assumed to be true. Otherwise, the modern cryptosystem can be broken in polynomial time. This type of hardness results hold for any algorithms.
- For orthonormal classes, we show that the gradient variance (wrt the target function) of is exponentially small. Hence, *gradient-based* algorithms are unlikely to succeed. This observation hold for any parametric model as long as they satisfy some smooth condition.
- Typical examples include the parity function and the cosine neuron: $f(x) = \cos(w^T x)$. The Barron norms of these functions are not greater than $O(d^2)$.

Summary

- Learning a subset of two-layer neural networks, whose path norms are bounded by $\text{poly}(d)$, can be reduced to certain classical hard problems, whose hardness is assumed to be true. Otherwise, the modern cryptosystem can be broken in polynomial time. This type of hardness results hold for any algorithms.
- For orthonormal classes, we show that the gradient variance (wrt the target function) of is exponentially small. Hence, *gradient-based* algorithms are unlikely to succeed. This observation hold for any parametric model as long as they satisfy some smooth condition.
- Typical examples include the parity function and the cosine neuron: $f(x) = \cos(w^T x)$. The Barron norms of these functions are not greater than $O(d^2)$.
- These hardness results suggest that the Barron space is very likely too large to study the training of two-layer neural networks.

Convergence of GD in the kernel regime

Setup

- Consider the two-layer neural network (2LNN):

$$f_m(x; \theta) = \sum_{j=1}^m a_j \sigma(b_j^T x), \quad (0.18)$$

where $\theta = (a, B)$ be the parameters, and $\sigma(z) = \max(0, z)$. The results can be extended to general Lipschitz activation functions with small modifications.

- The empirical risk with the square loss is given by

$$\hat{\mathcal{R}}_n(\theta) = \frac{1}{2n} \sum_{i=1}^n (f_m(x_i; \theta) - y_i)^2. \quad (0.19)$$

- Let $e_i = f_m(x_i; \theta) - y_i$. The GD flow is given by

$$\begin{aligned} \dot{a}_j &= - \sum_{i=1}^n e_i \sigma(b_j^T x_i) \\ \dot{b}_j &= - \sum_{i=1}^n e_i a_j \sigma'(b_j^T x_i) x_i. \end{aligned} \quad (0.20)$$

- Let $\pi_0 = \text{Unif}(\mathbb{S}^{d-1})$. We will mainly focus on the initialization:

$$a_j = 0, \quad b_j \sim \pi_0, \quad \text{for } j = 1, \dots, m. \quad (0.21)$$

Convergence results

Define an associate kernel

$$k(x, x') = \mathbb{E}_{b \sim \pi_0} [\sigma(b^T x) \sigma(b^T x')].$$

Let $K = (k(x_i, x_j)) \in \mathbb{R}^{n \times n}$ be the kernel matrix, and $\lambda_n(K)$ be the smallest eigenvalue of K .

Theorem 10

Let $\theta(t)$ be the GD solution at time t . For any $\delta \in (0, 1)$, assume that $m \geq \frac{10n \log(2n^2/\delta)}{\lambda_n^2(K)}$. Then, with probability $1 - \delta$ over the initialization, we have

$$\hat{\mathcal{R}}_n(\theta(t)) \leq e^{-m\lambda_n t} \hat{\mathcal{R}}_n(\theta_0).$$

A general observation

The proof relies on the following observation.

- The GD flow can be written as

$$\dot{\theta}(t) = -\frac{1}{n} \sum_{i=1}^n (f(x_i; \theta(t)) - y_i) \nabla_{\theta} f(x_i; \theta(t))$$

- Let $e_i(t) = f(x_i; \theta(t)) - y_i$. Then, we have

$$\frac{de_i(t)}{dt} = \langle \nabla f(x_i; \theta(t)), \dot{\theta}(t) \rangle = -\sum_{i'=1}^n \langle \nabla f(x_i; \theta(t)), \nabla f(x_{i'}; \theta(t)) \rangle e_{i'}(t).$$

- Let $G = (G_{i,i'}) \in \mathbb{R}^{n \times n}$ with $G_{i,i'} = \langle \nabla f(x_i; \theta), \nabla f(x_{i'}; \theta) \rangle$ be the Gram matrix, and $e = (e_1, \dots, e_n)^T \in \mathbb{R}^n$. Then,

$$\frac{de(t)}{dt} = -G(\theta(t))e(t). \quad (0.22)$$

- If $\lambda_n(G(\theta_t))$ is bounded away from zero for any $t \geq 0$, then the empirical risk converges to zero exponentially fast, since

$$\frac{d\|e(t)\|^2}{dt} = -2e(t)^T G(\theta_t)e(t) \leq -2\lambda_n(G(\theta_t))\|e(t)\|^2. \quad (0.23)$$

The Gram matrix for a 2LNN

- In this case, $G(\theta) = m\hat{K}(\theta)$

$$\hat{K}_{i,i'}(\theta) = \frac{1}{m} \sum_{j=1}^m \sigma(b_j^T x_i) \sigma(b_j^T x_{i'}) + a_j^2 \sigma'(b_j^T x_i) \sigma'(b_j^T x_{i'}) x_i^T x_{i'}.$$

- As $m \rightarrow \infty$, we have

$$\hat{K}_{i,i'} \rightarrow K_{i,i'} = k(x_i, x_{i'}).$$

where

$$k(x, x') = \mathbb{E}_{a,b}[\sigma(b^T x) \sigma(b^T x') + a^2 \sigma'(b^T x) \sigma'(b^T x') x^T x']. \quad (0.24)$$

- For the initialization considered here,

$$k(x, x') = \mathbb{E}_{b \sim \pi_0}[\sigma(b^T x) \sigma(b^T x')]. \quad (0.25)$$

Here, only the gradients wrt a contribute to the kernel.

Positivity of the Gram matrix at initialization

Lemma 11

Assume $\lambda_n(K) > 0$. For any $\delta \in (0, 1)$, if $m \geq \frac{\log(n^2/\delta)}{2\lambda_n^2(K)}$, with probability $1 - \delta$ over the random initialization, we have

$$\lambda_n(G) \geq \frac{m}{2} \lambda_n(K).$$

Remark:

- The condition: $\lambda_n(K) > 0$ does not allow two samples x_i and x_j to align with each other for $i \neq j$.
- If $\{x_i\}_{i=1}^n$ are independently drawn from $\text{Unif}(\mathbb{S}^{d-1})$, [Braun, 2006] proved that with high probability, $\lambda_n(K) > n\lambda_n/2$, where λ_n is the n -th largest eigenvalue of the kernel function $k(\cdot, \cdot)$.
- For the ReLU activation function, one can show that $\lambda_n \geq \frac{C_d}{n^{1+1/d}}$. So, $\lambda_n(K) \geq C_d/n^{1/d}$ (see the appendix of (Ma et al, MSML2020)).
- We will leave $\lambda_n(K) > 0$ as a basic assumption.

Positivity of the Gram matrix at initialization

Proof:

- By Hoeffding's inequality, we have for any $i, j \in [n]$

$$\mathbb{P}\{|\hat{K}_{i,j} - K_{i,j}| \geq \varepsilon\} = \mathbb{P}\left\{\left|\frac{1}{m} \sum_{s=1}^m \sigma(b_s^T x_i) \sigma(b_s^T x_j) - \mathbb{E}[\sigma(b_s^T x_i) \sigma(b_s^T x_j)]\right| \geq \varepsilon\right\} \leq e^{-2m\varepsilon^2}.$$

- Taking the union bound leads to

$$\mathbb{P}\{\|\hat{K} - K\|_F \leq \varepsilon\} \geq 1 - \sum_{i,j=1}^n \mathbb{P}\{|\hat{K}_{i,j} - K_{i,j}| \geq \varepsilon\} \geq 1 - n^2 e^{-2m\varepsilon^2}.$$

- Using the Weyl's inequality, we have

$$\lambda_n(\hat{K}) \geq \lambda_n(K) - \|\hat{K} - K\|_F \geq \lambda_n(K) - \varepsilon.$$

- Take $\varepsilon = \lambda_n(K)/2$ and let the failure prob. $n^2 e^{-2m\varepsilon^2} \leq \delta$. This leads to $m \geq \frac{\log(n^2/\delta)}{2\lambda_n^2(K)}$.

Gradient descent near the initialization

- Define a neighbor of the initialization by

$$\mathcal{I}(\theta_0) := \left\{ \theta : \|\hat{K}(\theta) - \hat{K}(\theta_0)\|_F \leq \frac{\lambda_n(K)}{4} \right\}$$

- Using Lemma 11, for any $\delta \in (0, 1)$ with probability $1 - \delta$, we have for any $\theta \in \mathcal{I}(\theta_0)$ that

$$\lambda_n(\hat{K}(\theta)) \geq \lambda_n(\hat{K}(\theta_0)) - \|\hat{K}(\theta_0) - \hat{K}(\theta)\|_F \geq \frac{\lambda_n(K)}{2} - \frac{\lambda_n(K)}{4} = \frac{\lambda_n(K)}{4}.$$

Lemma 12

Let $t_0 = \inf \{t : \theta(t) \notin \mathcal{I}(\theta_0)\}$. For any $\delta \in (0, 1)$, assume $m \geq \frac{\log(n^2/\delta)}{2\lambda_n^2(K)}$. For any $t \in [0, t_0]$,

$$\hat{\mathcal{R}}_n(\theta(t)) \leq e^{-\frac{m\lambda_n(K)}{2}t} \hat{\mathcal{R}}_n(\theta_0).$$

Proof:

$$\frac{d}{dt} \hat{\mathcal{R}}_n(\theta(t)) = \frac{1}{2n} \frac{d\|e(t)\|^2}{dt} = \frac{-m}{n} e^T \hat{K} e \leq \frac{-m}{n} \frac{\lambda_n(K)}{4} \|e(t)\|^2 = \frac{-m\lambda_n(K)}{2} \hat{\mathcal{R}}_n(\theta_t)$$

Long-time convergence of GD

Proof of Theorem 10:

- We only need to prove that $t_0 = \infty$. Otherwise, assume that $t_0 < \infty$.
- First, the empirical risk is *smooth* in the sense that

$$\|\nabla \hat{\mathcal{R}}_n(\theta)\|^2 \leq \|\theta\|^2 \hat{\mathcal{R}}_n(\theta).$$

- Then,

$$\begin{aligned} \|\theta(t) - \theta_0\| &\leq \int_0^{t_0} \|\nabla \hat{\mathcal{R}}_n(\theta(t))\| dt \leq \max_{t \in [0, t_0]} \|\theta(t)\| \int_0^{t_0} \sqrt{\hat{\mathcal{R}}_n(\theta)} dt \\ &\leq \max_{t \in [0, t_0]} \|\theta(t)\| \int_0^{t_0} e^{-\frac{m\lambda_n(K)}{4}t} \sqrt{\hat{\mathcal{R}}_n(\theta_0)} dt \lesssim \frac{\max_{t \in [0, t_0]} \|\theta(t)\|}{m\lambda_n(K)}. \end{aligned}$$

Let $\gamma = \max_{t \in [0, t_0]} \|\theta(t) - \theta_0\|$. Using the fact that $\|\theta_0\| = \sqrt{\sum_{s=1}^m \|b_s(0)\|^2} = \sqrt{m}$, we have

$$\gamma \lesssim \frac{\gamma + \sqrt{m}}{m\lambda_n(K)},$$

which leads to

$$\gamma \lesssim \frac{1}{\sqrt{m}\lambda_n(K)}.$$

Long-time convergence of GD (Cont'd)

Proof of Theorem 10:

- Since σ is 1-Lipschitz continuous, we have for any $t \in [0, t_0]$,

$$\begin{aligned}\|\hat{K}(\theta(t)) - K(\theta_0)\|_F^2 &= \sum_{i,j} \left| \frac{1}{m} \sum_{s=1}^m \sigma(b_s(t)^T x_i) \sigma(b_s(t)^T x_j) - \frac{1}{m} \sum_{s=1}^m \sigma(b_s(0)^T x_i) \sigma(b_s(0)^T x_j) \right|^2 \\ &\lesssim \frac{n^2}{m^2} (\|\theta(t) - \theta_0\| + \|\theta(t) - \theta_0\|^2) \\ &\lesssim \frac{n^2}{m^2} (\gamma + \gamma^2).\end{aligned}$$

- By the assumption, $\gamma \leq 1$. Hence, $m \geq 20n/\lambda_n(K)$ leads to

$$\|\hat{K}(\theta(t)) - K(\theta_0)\|_F \leq \frac{\lambda_n(K)}{8},$$

which contradicts the definition of t_0 . Therefore, $t_0 = \infty$.

Remarks

- In the above analysis, the main ingredient is the positivity of the Gram matrix G , which relies on the positivity of the **tangent kernel**:

$$k(x, x') = \lim_{m \rightarrow \infty} m^{-\alpha} \langle \nabla f(x; \theta), \nabla f(x'; \theta) \rangle,$$

where α is a specific factor related to the initialization such that the limit exists.

- The key observation is that $b_j(t) - b_j(0) \sim \frac{1}{m}$. The parameters of the convergent solution is close to the initialization.
- The results can be extended to general initializations. For instance, consider the balanced initialization: $a_j \sim \mathcal{N}(0, 1/m)$, $b_j \sim \mathcal{N}(0, I_d/(md))$, for which the Gram matrix

$$G_{i,i'} = \sum_{j=1}^m \sigma(b_j^T x_i) \sigma(b_j^T x_{i'}) + a_j^2 \sigma'(b_j^T x_i) \sigma'(b_j^T x_{i'}) x_i^T x_{i'}$$

$$\rightarrow k(x_i, x_{i'}) := \mathbb{E}_{b \sim \mathcal{N}(0, I_d/d)} [\sigma(b^T x_i) \sigma(b^T x_{i'}) + \sigma'(b^T x_i) \sigma'(b^T x_{i'}) x_i^T x_{i'}] \text{ as } m \rightarrow \infty.$$

We only need to show that smallest eigenvalue of the kernel matrix:

$K = (k(x_i, x_{i'})) \in \mathbb{R}^{n \times n}$ is away from zero.

Characterization of the whole GD trajectory

The following theorem concerns the function class that the GD solutions can represent. Let $f_m(x; a, B_0) = \sum_{j=1}^m a_j \sigma(b_j(0)^T x)$ be the random feature model (RFM).

Theorem 13 (E, Ma, Wu, 2019)

Let $\theta_t = \{a(t), B(t)\}$ be the GD solution at time t , and $\tilde{a}(t)$ be the GD solution of RFM with zero initialization. For any $\delta \in (0, 1)$, assume that $m \gtrsim \frac{n^4}{\lambda_n^2(K)} \log(\frac{n^2}{\delta})$. Then, with probability $1 - \delta$ over the random initialization, we have

$$\sup_{t \in [0, \infty], x \in \mathbb{S}^{d-1}} |f_m(x; a(t), B(t)) - f_m(x; \tilde{a}(t); B_0)| \leq \frac{1 + \sqrt{\log(1/\delta)}}{\lambda_n(K) \sqrt{m}}.$$

Remark:

- The theorem implies that the GD trajectory of a wide 2LNN is **uniformly** close to that of the associate RFM.
- The result is implicit in the proof of convergence result: $\theta(t) - \theta_0 \ll 1$.
- **Time-scale separation:** To fit n labels, we only need $a_j(t) = O(\text{poly}(n)/m)$. Then, $\dot{a}_j(t) \sim O(\|b_j\|) = O(1)$ and $\dot{b}_j(t) \sim O(|a_j|) = O(\text{poly}(n)/m)$. Hence, $b_j(t)$ is essentially frozen when $m \rightarrow \infty$.

Proof sketch

$$\begin{aligned} |f_m(x; a(t), B(t)) - f_m(x; a(t), B_0)| &\leq \sum_{s=1}^m a_s(t) |\sigma(b_s(t)^T x) - \sigma(b_s(0)^T x)| \\ &\leq \sum_{s=1}^m a_s(t) \|b_s(t) - b_s(0)\| \\ &\leq \frac{1}{2} \sum_{s=1}^m (a_s^2(t) + \|b_s(t) - b_s(0)\|^2) \\ &= \frac{1}{2} \|\theta(t) - \theta_0\|^2 \lesssim \frac{1}{m\lambda_n^2(K)} \end{aligned}$$

The closeness of $a(t)$ and $\tilde{a}(t)$ is a consequence of the closeness of $B(t)$ and B_0 . The proof is lengthy but straightforward.

Representer theorem of GD solutions

Lemma 14

Let $k_m(x, x') = \frac{1}{m} \sum_{j=1}^m \sigma(b_j(0)^T x) \sigma(b_j(0)^T x')$. Then, there exist $w_1(t), \dots, w_n(t)$ such that the GD solution of RFM with zero initialization can be written as

$$f_m(x; \tilde{a}(t), B_0) = \sum_{i=1}^n w_i(t) k_m(x_i, x).$$

Proof: Note that

$$\frac{d\tilde{a}_j(t)}{dt} = - \sum_{i=1}^n e_i(t) \sigma(b_j(0)^T x_i).$$

Hence,

$$\begin{aligned} f(x; \tilde{a}(t), B_0) &= \sum_{j=1}^m a_j(t) \sigma(b_j(0)^T x) = - \sum_{j=1}^m \left(\int_0^t \sum_{i=1}^n e_i(t') dt' \sigma(b_j(0)^T x_i) \right) \sigma(b_j(0)^T x) \\ &= \sum_{i=1}^n w_i(t) k_m(x_i, x), \end{aligned} \tag{0.26}$$

where $w_i(t) = -m \int_0^t e_i(t') dt'$.

Curse of dimensionality

- Lemma 14 shows that the GD solutions of RFM always lie in the span of $\{k_m(x_i, \cdot)\}_{i=1}^n$ no matter how big m is.
- Recall the definition of Barron space

$$\|f\|_{\mathcal{B}} = \inf_{f(x) = \mathbb{E}_{(a,b) \sim \rho} [a\sigma(b^T x)]} \mathbb{E}[|a|\|b\|].$$

- With explicit regularization, the generalization error for learning Barron functions obeys the Monte-Carlo rate: $O(1/m + 1/\sqrt{n})$.
- However, [Barron 1993] shows that for any $h_1, \dots, h_n \in L^2(X)$

$$\sup_{\|f\|_{\mathcal{B}} \leq 1} \inf_{h \in \text{span}(h_1, \dots, h_n)} \|f - h\|_{L^2(X)} \geq \frac{C}{dn^{1/d}}.$$

Therefore, GD solutions of 2LNN suffer from the curse of dimensionality.

- Note that the uniform closeness implies that any early stopping “regularization” cannot cure the CoD.

Compare the NN and RFM under GD dynamics

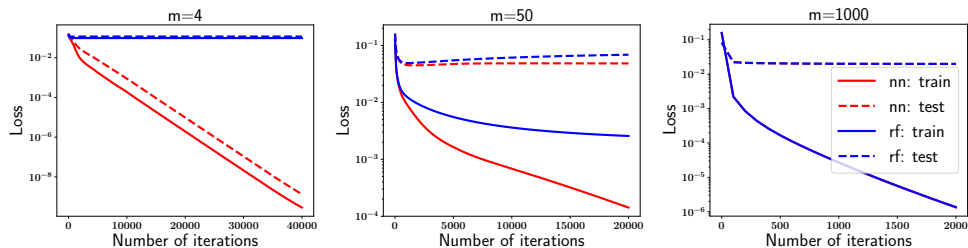


Figure 2: GD dynamics for fitting a single neuron $f^*(x) = \sigma(x_1)$ where $x \in \mathbb{S}^{d-1}$. Here $d = 20, n = 50$. **Left:** $m = 4$; **Middle:** $m = 50$; **Right:** $m = 1000$.

Comparison between the implicit and explicit regularization

Consider the explicit regularization:

$$\min_{\theta} \hat{\mathcal{R}}_n(\theta) + \lambda \sqrt{\frac{\log(d)}{n}} \sum_{j=1}^m |a_j| \|b_j\|_2.$$

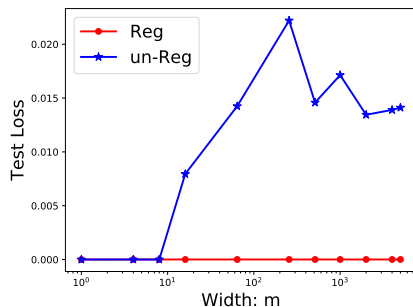


Figure 3: Fitting a single neuron $f^*(x) = \sigma(x_1)$ where $x \in \mathbb{S}^{d-1}$. Here $d = 20, n = 50$. The nn-Reg solution is the GD solution without any regularization.

Lazy training

- In the previous analysis, the main insight is that for the highly over-parameterized setting, the perturbation satisfies $\|\theta - \theta_0\| \ll 1$.
- Why is a small deviation enough? Consider the expansion around the initialization:

$$f_m(x; \theta) = f_m(x; \theta_0) + \langle \theta - \theta_0, \nabla f_m(x; \theta_0) \rangle + o(\|\theta - \theta_0\|). \quad (0.27)$$

- Note that the each entry of $\theta - \theta_0$ is in the order of $O(1/m)$ is enough to ensure the change of output: $f_m(x; \theta) - f_m(x; \theta_0) \sim 1$. Meanwhile, $\|\theta - \theta_0\| = O(1/\sqrt{m}) \ll 1$.
- So essentially, only the linear part contributes to the final model. In our case, the linear part is a RFM.
- In the literature, training methods that essentially only explore the linear part of a nonlinear model to fit data are called **lazy training** (Chizat, Oyallon, Bach, 2018).

Neural tangent kernel

- In the lazy training regime, the model essentially performs kernel method with the kernel given by:

$$k_m(x, x') = \langle \nabla f_m(x; \theta_0), \nabla f_m(x'; \theta_0) \rangle,$$

which is called tangent kernel.

- Large width limit: For neural network models $f_m(\cdot; \theta)$, k_m often has a limit with a proper rescaling:

$$k(x, x') = \lim_{m \rightarrow \infty} m^{-\alpha} k_m(x, x').$$

$k(\cdot, \cdot)$ is called the **neural tangent kernel** (NTK) (Jacot, Gabriel and Hongler, 2018).

Multi-layer fully-connected networks

The observation that GD only performs lazy training can be extended to general wide neural networks. The proof is similar to the case of 2LNN and can be summarized as follows.

- Recall that

$$\frac{d\|e(t)\|^2}{dt} = -2e(t)^T G(\theta(t))e(t). \quad (0.28)$$

- First show that if m is sufficiently large, at initialization $\lambda_n(G(\theta_0)) \geq m\lambda_n(K)$, where K is the kernel matrix of NTK. Assume that $\lambda_n(K) > 0$ ([Justify it](#)).
- Let $I(\theta_0)$ be the ball around the initialization such that the $\lambda_n(G(\theta)) \geq m\lambda_n(K)/2$. Let t_0 be the time that $\theta(t)$ first leaves the ball. Then, for any $t \in [0, t_0]$, we have $\hat{R}_n(\theta(t)) \leq e^{-cm\lambda_n(K)t} \hat{R}_n(\theta_0)$ for a constant $c > 0$.
- The combination of exponential convergence and continuity of \hat{R}_n implies

$$\|\theta_t - \theta_0\| \leq \int_0^{t_0} \|\nabla \hat{R}_n(\theta'_t)\| dt' \leq \int_0^{t_0} C(\|\theta(t)\|) \sqrt{\hat{R}_n(\theta(t'))} dt' \leq \frac{\text{poly}(n, \lambda_n(K))}{m}.$$

- When m is sufficiently large, we must have $\theta_t \in \mathcal{I}(\theta_0)$ for any $t \geq 0$.

We refer to ([Arora et al, 2019](#)) for a detailed proof for multi-layer fully-connected networks.

Summary

- Under conventional setting, neural networks trained by GD converges to kernel method. Moreover, the convergence is uniform in time. It means only when f^* lies in the appropriate RKHS, the GD solution can generalize well.
- What happens when the network is less over-parameterized?
- Can we still learn larger class of target functions in the over-parameterized regime?