

On the Hardness of Learning Two-layer Neural Networks

Instructor: Weinan E

Mathematical Introduction to Machine Learning
MAT 490/APC 490

Princeton University, Spring 2021

Motivation

- For two-layer neural networks, we have defined the Barron spaces. Both the approximation and estimation error for target functions in these spaces obey the Monte-Carlo rate, which is free of the curse of dimensionality (CoD). A nature questions is then: Do there exist algorithms can learn these functions efficiently?

Motivation

- For two-layer neural networks, we have defined the Barron spaces. Both the approximation and estimation error for target functions in these spaces obey the Monte-Carlo rate, which is free of the curse of dimensionality (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 for returning a solution \hat{f} such that $\|\hat{f} - f^*\| \leq \varepsilon$ satisfies:

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

Motivation

- For two-layer neural networks, we have defined the Barron spaces. Both the approximation and estimation error for target functions in these spaces obey the Monte-Carlo rate, which is free of the curse of dimensionality (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 for returning a solution \hat{f} such that $\|\hat{f} - f^*\| \leq \varepsilon$ satisfies:

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

- Note that in the highly over-parameterized setting, GD can find a global minimum of the empirical risk exponentially fast. However, this solution is essentially a kernel predictor, whose generalization error suffers from CoD. Hence, the total time complexity still suffers from CoD.

Motivation

- For two-layer neural networks, we have defined the Barron spaces. Both the approximation and estimation error for target functions in these spaces obey the Monte-Carlo rate, which is free of the curse of dimensionality (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 for returning a solution \hat{f} such that $\|\hat{f} - f^*\| \leq \varepsilon$ satisfies:

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

- Note that in the highly over-parameterized setting, GD can find a global minimum of the empirical risk exponentially fast. However, this solution is essentially a kernel predictor, whose generalization error suffers from CoD. Hence, the total time complexity still suffers from CoD.
- In this lecture, 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 $\sigma_{0,1}$ be the step function, i.e., $\sigma_{0,1}(t) = 1(t \geq 0)$.

Learning intersections of halfspaces

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

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

Theorem 1 (Theorem 1.2, [Kalai, Klivans, 2008](#))

Let $\mathcal{H} = \{\mathbf{x} \mapsto \sigma_{0,1}(\mathbf{w}^T \mathbf{x} - b - 1/2) : b \in \mathbb{N}, \mathbf{w} \in \mathbb{N}^d, |b| + \|\mathbf{w}\|_1 \leq \text{poly}(d)\}$. Define

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

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 will show two-layer neural networks can simulate the functions in \mathcal{H}_K .

Hardness of learning two-layer ReLU networks

Theorem 2 (Livni, et al, 2014)

Let $\mathcal{X} = \{-1, 1\}^d$, and \mathcal{G} be the class of functions with the Barron norm bounded by $\text{poly}(d)$. 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:

$$\sigma_{0,1}(t) = \lim_{a \rightarrow \infty} \text{ReLU}(at) - \text{ReLU}(at - 1).$$

Hardness of learning two-layer neural networks

Proof:

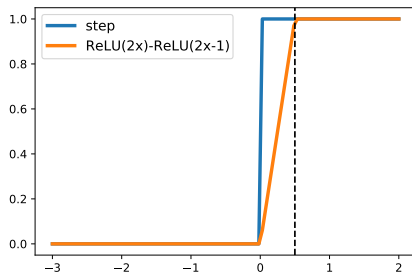
- Let $c(\mathbf{x}) = \mathbf{w}^T \mathbf{x} - b - 1/2$. Since $\mathbf{w} \in \mathbb{N}^d, \mathbf{x} \in \{-1, 1\}^d, b \in \mathbb{N}$, we have $|c(\mathbf{x})| \geq 1/2$. Assume $\|\mathbf{w}\|_1 + |b| \leq \text{poly}(d)$.
- Consider k hyperplanes $\{c_i\}_{i=1}^k$. Let $h_i(\mathbf{x}) = \sigma_{0,1}(c_i(\mathbf{x})) \in \mathcal{H}$.
- Let

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

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

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

- The blue part is equal to $\sigma_{0,1}(c_i(\mathbf{x}))$ due to $\sigma_{0,1}(z) = \text{ReLU}(2z) - \text{ReLU}(2z - 1)$ for $|z| \geq 1/2$.



- We can verify that

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

- 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 $\langle f_1, f_2 \rangle = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[f_1(\mathbf{x})f_2(\mathbf{x})]$.

Definition 3 (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 $\langle f_1, f_2 \rangle = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[f_1(\mathbf{x})f_2(\mathbf{x})]$.

Definition 3 (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 $f_i, f_j \in \mathcal{F}$.

Definition 4 (Statistical query (SQ) dimension)

Let \mathcal{F} be a function class with $\|f\| = 1$ for any $f \in \mathcal{F}$. The SQ dimension of \mathcal{F} is the largest n such that there is a n -sized subset $\{f_i\}_{i=1}^n \subset \mathcal{F}$ such that $\langle f_i, f_j \rangle \leq \frac{1}{n^2}$ for any $i \neq j$ and $i, j \in [n]$.

The SQ dimension characterizes the number of “nearly” orthonormal functions 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 $\langle f_1, f_2 \rangle = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[f_1(\mathbf{x})f_2(\mathbf{x})]$.

Definition 3 (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 $f_i, f_j \in \mathcal{F}$.

Definition 4 (Statistical query (SQ) dimension)

Let \mathcal{F} be a function class with $\|f\| = 1$ for any $f \in \mathcal{F}$. The SQ dimension of \mathcal{F} is the largest n such that there is a n -sized subset $\{f_i\}_{i=1}^n \subset \mathcal{F}$ such that $\langle f_i, f_j \rangle \leq \frac{1}{n^2}$ for any $i \neq j$ and $i, j \in [n]$.

The SQ dimension characterizes the number of “nearly” orthonormal functions in \mathcal{F} .

- Let $\mathcal{B}_d = \{f_m(\cdot; \theta) : \|\theta\|_{\mathcal{P}} \leq Cd^2, m \in \mathbb{N}\}$. Here, f_m is a two-layer neural network and $C > 0$ is sufficiently large. We will show that the SQ dimension of \mathcal{B}_d is exponentially in d .
- We will show that if the SQ dimension of a function class \mathcal{F} is exponentially large, the learning of \mathcal{F} will be hard in some sense.

Parity functions

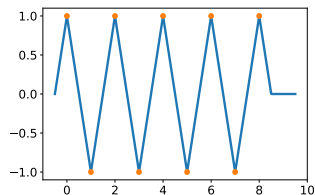
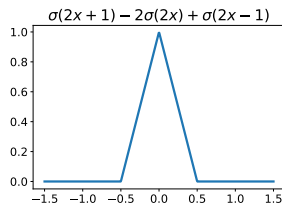
- $\mathcal{F}_1 = \{f_{\mathbf{v}}(\mathbf{x}) = (-1)^{\langle \mathbf{v}, \mathbf{x} \rangle} : \mathbf{v} \in \{0, 1\}^d\}$, where $\mathbf{x} \in \{0, 1\}^d$.
- Consider $\mathcal{D} = \text{Unif}(\{0, 1\}^d)$. Then, we have

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

$$= \prod_{i=1}^d \mathbb{E}_{x_i}[(-1)^{(v_i + v'_i)x_i}] = \delta_{\mathbf{v}, \mathbf{v}'} \quad (2)$$

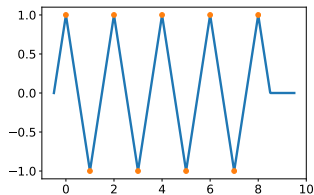
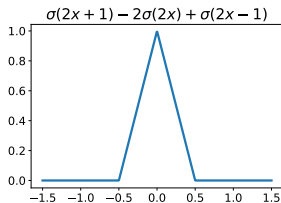
Hence, \mathcal{F}_1 is an orthonormal class with $|\mathcal{F}_1| = 2^d$.

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 $\mathbf{v}^T \mathbf{x} \in \{0, 1, \dots, d\}$. Let σ be the ReLU function. Then,

$$(-1)^{\mathbf{v}^T \mathbf{x}} = \sigma_{tri}(\mathbf{v}^T \mathbf{x}) = \sum_{i=0}^d (-1)^i \sigma_{hat}(\mathbf{v}^T \mathbf{x} - i) \quad (3)$$

$$= \sum_{i=0}^d (-1)^i (\sigma(2\mathbf{v}^T \mathbf{x} - 2i + 1) - 2\sigma(2\mathbf{v}^T \mathbf{x} - 2i) + \sigma(2\mathbf{v}^T \mathbf{x} - 2i - 1)) \quad (4)$$

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

The cosine neuron

Cosine neuron: Let $\mathcal{D} = \mathcal{N}(0, I_d)$. Define

$$\mathcal{F}_2 := \left\{ f_{\mathbf{w}}(\mathbf{x}) = \cos(2\pi \mathbf{w}^T \mathbf{x}) : \|\mathbf{w}\| = \sqrt{d} \right\}.$$

Then the SQ dimension of \mathcal{F}_2 is larger than $C_1 e^{C_2 \sqrt{d}}$, where C_1, C_2 are absolute positive constants.

The cosine neuron (Cont'd)

Proof:

- Let $\mathbf{w}_i \stackrel{iid}{\sim} \mathbb{S}^{d-1}$ with $i = 1, \dots, n$. Then, taking the union bound gives us

$$\mathbb{P} \left\{ \max_{i \neq j} |\mathbf{w}_i^T \mathbf{w}_j| \geq t \right\} \leq 2n^2 e^{-C_1 d t^2}.$$

Let $t = 1/d^{1/4}$ and $2n^2 e^{-C_1 t^2} \leq 1/2$. We have $n = e^{-C_2 d^{1/2}}/2$.

- By the rescaling to $\sqrt{d}\mathbb{S}^{d-1}$, it is implied that there exist $n = 0.5e^{-C_2 d^{1/2}}$ points $\mathbf{w}_1, \dots, \mathbf{w}_n$ such that

$$\|\mathbf{w}_i\|^2 = d, \quad |\langle \mathbf{w}_i, \mathbf{w}_j \rangle| \leq \sqrt{d}\sqrt{d} \frac{1}{d^{1/4}} = d^{3/4}, \quad \forall i \neq j.$$

- Denote by $\cos_{\mathbf{w}}(\mathbf{x}) = \cos(2\pi \mathbf{w}^T \mathbf{x})$ and φ^2 the density function of the standard Gaussian. Then, we have for $i \neq j$,

$$\begin{aligned}
 \langle \cos_{\mathbf{w}_i}, \cos_{\mathbf{w}_j} \rangle &= \int \varphi^2(\mathbf{x}) \cos(\mathbf{w}_i^T \mathbf{x}) \cos(\mathbf{w}_j^T \mathbf{x}) d\mathbf{x} = \langle \cos_{\mathbf{w}_i} \varphi, \cos_{\mathbf{w}_j} \varphi \rangle_{L^2(\mathbb{R}^d)} \\
 &= \langle \widehat{\cos_{\mathbf{w}_i}} * \hat{\varphi}, \widehat{\cos_{\mathbf{w}_j}} * \hat{\varphi} \rangle_{L^2(\mathbb{R})} \\
 &= \frac{1}{4} \int_{\mathbb{R}^d} (\hat{\varphi}(\boldsymbol{\xi} - \mathbf{w}_i) + \hat{\varphi}(\boldsymbol{\xi} + \mathbf{w}_i)) (\hat{\varphi}(\boldsymbol{\xi} - \mathbf{w}_j) + \hat{\varphi}(\boldsymbol{\xi} + \mathbf{w}_j)) d\boldsymbol{\xi}.
 \end{aligned}$$

Here, $\hat{\varphi}(\boldsymbol{\xi}) = \frac{(4\pi)^{d/2}}{(2\pi)^{d/4}} e^{-4\pi^2 \|\boldsymbol{\xi}\|^2}$ is a Gaussian-like function. Since,

$$\|\mathbf{w}_i - \mathbf{w}_j\|^2 = 2d - 2\langle \mathbf{w}_i, \mathbf{w}_j \rangle \geq 2d - 2d^{3/4} \geq d/2 \text{ for sufficiently large } d.$$

- Denote by $\cos_{\mathbf{w}}(\mathbf{x}) = \cos(2\pi \mathbf{w}^T \mathbf{x})$ and φ^2 the density function of the standard Gaussian. Then, we have for $i \neq j$,

$$\begin{aligned}
\langle \cos_{\mathbf{w}_i}, \cos_{\mathbf{w}_j} \rangle &= \int \varphi^2(\mathbf{x}) \cos(\mathbf{w}_i^T \mathbf{x}) \cos(\mathbf{w}_j^T \mathbf{x}) d\mathbf{x} = \langle \cos_{\mathbf{w}_i} \varphi, \cos_{\mathbf{w}_j} \varphi \rangle_{L^2(\mathbb{R}^d)} \\
&= \langle \widehat{\cos_{\mathbf{w}_i}} * \hat{\varphi}, \widehat{\cos_{\mathbf{w}_j}} * \hat{\varphi} \rangle_{L^2(\mathbb{R})} \\
&= \frac{1}{4} \int_{\mathbb{R}^d} (\hat{\varphi}(\boldsymbol{\xi} - \mathbf{w}_i) + \hat{\varphi}(\boldsymbol{\xi} + \mathbf{w}_i))(\hat{\varphi}(\boldsymbol{\xi} - \mathbf{w}_j) + \hat{\varphi}(\boldsymbol{\xi} + \mathbf{w}_j)) d\boldsymbol{\xi}.
\end{aligned}$$

Here, $\hat{\varphi}(\boldsymbol{\xi}) = \frac{(4\pi)^{d/2}}{(2\pi)^{d/4}} e^{-4\pi^2 \|\boldsymbol{\xi}\|^2}$ is a Gaussian-like function. Since,

$$\|\mathbf{w}_i - \mathbf{w}_j\|^2 = 2d - 2\langle \mathbf{w}_i, \mathbf{w}_j \rangle \geq 2d - 2d^{3/4} \geq d/2 \text{ for sufficiently large } d.$$

- Hence, there exist a constant $C_3 > 0$ such that

$$\langle \cos_{\mathbf{w}_i}, \cos_{\mathbf{w}_j} \rangle \leq \sup_{\|\mathbf{w} - \mathbf{v}\|^2 \geq C_3 d} \int \hat{\varphi}(\boldsymbol{\xi} + \mathbf{w}) \hat{\varphi}(\boldsymbol{\xi} + \mathbf{v}) d\boldsymbol{\xi} \quad (5)$$

$$= e^{-2\pi^2 \|\mathbf{w} - \mathbf{v}\|^2} \leq e^{-2\pi^2 C_3 d} \leq \frac{1}{n^2}. \quad (6)$$

Here, $n = 0.5e^{-C_2 \sqrt{d}}$. Therefore, the SQ dimension of \mathcal{F}_2 is at least exponentially in \sqrt{d} .

Approximating the cosine neuron with two-layer neural networks

- Note that Barron function is defined on a compact domain. We can consider $\mathcal{D} = \text{Unif}(\sqrt{d}\mathbb{S}^{d-1})$. This does not change the analysis too much since when $d \gg 1$, $\mathcal{N}(0, I_d)$ is similar to $\text{Unif}(\sqrt{d}\mathbb{S}^{d-1})$.

Approximating the cosine neuron with two-layer neural networks

- Note that Barron function is defined on a compact domain. We can consider $\mathcal{D} = \text{Unif}(\sqrt{d}\mathbb{S}^{d-1})$. This does not change the analysis too much since when $d \gg 1$, $\mathcal{N}(0, I_d)$ is similar to $\text{Unif}(\sqrt{d}\mathbb{S}^{d-1})$.
- Through rescaling, we can consider the target function $f^*(\mathbf{x}) = \cos(\sqrt{d}\mathbf{w}^T \mathbf{x})$ and $\mathbf{x} \in [-1, 1]^d$. By the spectral characterization of Barron functions, we have

$$\|f^*\|_{\mathcal{B}} \leq 2 \int_{\mathbb{R}^d} \|\boldsymbol{\xi}\|_1^2 |\mathcal{F}(f^*)(\boldsymbol{\xi})| d\boldsymbol{\xi} + 2|f^*(0)| + 2\|\nabla f(0)\|_1 \leq 2\|\sqrt{d}\mathbf{w}\|_1^2 + 2 \quad (7)$$

$$\leq Cd^3. \quad (8)$$

Remark

$\mathcal{F}_\phi := \{\phi(\mathbf{w}^T \mathbf{x}) : \|\mathbf{w}\| = \sqrt{d}\}$. $\mathcal{D} = \mathcal{N}(0, I_d)$. [Shamir, 2017] shows that as long as ϕ is periodic, under some mild condition, we have the SQ dimension of \mathcal{F}_ϕ is exponentially in d .

Gradients for an orthonormal class

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

Theorem 5

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

$$\mathbb{E}_{f \sim P}[(\nabla R^f(\theta) - g(\theta))^2] \leq \frac{\mathbb{E}_{\mathbf{x}}[\|\nabla_{\theta} h(\mathbf{x}; \theta)\|^2]}{|\mathcal{F}|}. \quad (9)$$

Gradients for an orthonormal class

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

Theorem 5

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

$$\mathbb{E}_{f \sim P}[(\nabla R^f(\theta) - g(\theta))^2] \leq \frac{\mathbb{E}_{\mathbf{x}}[\|\nabla_{\theta} h(\mathbf{x}; \theta)\|^2]}{|\mathcal{F}|}. \quad (9)$$

- 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 $R^f(\theta) = \mathbb{E}_{\mathbf{x}}[(h(\mathbf{x}; \theta) - f(\mathbf{x}))^2]$ the risk. Then, we have the following theorem.

Theorem 5

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

$$\mathbb{E}_{f \sim P}[(\nabla R^f(\theta) - g(\theta))^2] \leq \frac{\mathbb{E}_{\mathbf{x}}[\|\nabla_{\theta} h(\mathbf{x}; \theta)\|^2]}{|\mathcal{F}|}. \quad (9)$$

- 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

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

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

Proof:

- First, the gradient can be written as follows

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

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

- Hence,

$$\mathbb{E}_f[(\nabla_{\theta} R^f - g(\theta))^2] \leq \mathbb{E}_f[\langle f, \nabla_{\theta} h(\mathbf{x}; \theta) \rangle^2] \quad (10)$$

$$\leq \frac{1}{|\mathcal{F}|} \sum_f \langle f, \nabla_{\theta} h(\mathbf{x}; \theta) \rangle^2 \quad (11)$$

$$\leq \frac{\mathbb{E}_{\mathbf{x}}[\|\nabla_{\theta} h(\mathbf{x}; \theta)\|^2]}{|\mathcal{F}|}. \quad (12)$$

Extension to the nearly orthonormal classes

For the nearly orthonormal classes. Let n be the SQ dimension. Then, there exist f_1, \dots, f_n such that $\langle f_i, f_j \rangle \leq 1/n^2$ for $i \neq j$. Let $g = \nabla_{\theta} h(\cdot; \theta)$. WLOG, assume $\|g\| \leq 1$. Then,

$$0 \leq \|g - \sum_i \langle g, f_i \rangle f_i\|^2 \quad (13)$$

$$= \|g\|^2 - \sum_i \langle g, f_i \rangle^2 + \sum_{i \neq j} \langle g, f_i \rangle \langle g, f_j \rangle \langle f_i, f_j \rangle \quad (14)$$

$$\leq \|g\|^2 - \sum_i \langle g, f_i \rangle^2 + \frac{n^2}{n^2}. \quad (15)$$

Therefore, let P be the uniform distribution over f_1, \dots, f_n , we have

$$\mathbb{E}_{f \sim P}[(\nabla R^f(\theta) - g(\theta))^2] \leq \frac{1}{n} \sum_{i=1}^n \langle f_i, g \rangle^2 \leq \frac{2}{n}.$$

Therefore, we prove similar results for the nearly orthonormal classes, such as the cosine neuron class.

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

$$R^f(\theta) := \mathbb{E}_{\mathbf{x}}[\max(0, 1 - h(\mathbf{x}; \theta)f(\mathbf{x}))]. \quad (16)$$

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

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

Remark: the above assumption holds for parity functions.

Hardness of learning with GD

Theorem 6

Assume the model satisfies that $\sup_{\mathbf{x} \in X} |h(\mathbf{x}; \theta)| \leq 1$ and $\mathbb{E}_{\mathbf{x}}[\|\nabla_{\theta} h(\mathbf{x}; \theta_1) - \nabla_{\theta} h(\mathbf{x}; \theta_2)\|^2] \leq L\|\theta_1 - \theta_2\|^2$. 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

$$\mathbb{E}_f[\|\theta_t^f - \theta_0\|^2] \leq C_1(e^{\frac{C_2 t}{2^{d/2}}} - 1), \quad (17)$$

where C_1, C_2 only depend on L and θ_0 .

Hardness of learning with GD

Theorem 6

Assume the model satisfies that $\sup_{\mathbf{x} \in X} |h(\mathbf{x}; \theta)| \leq 1$ and $\mathbb{E}_{\mathbf{x}}[\|\nabla_{\theta} h(\mathbf{x}; \theta_1) - \nabla_{\theta} h(\mathbf{x}; \theta_2)\|^2] \leq L\|\theta_1 - \theta_2\|^2$. 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

$$\mathbb{E}_f[\|\theta_t^f - \theta_0\|^2] \leq C_1(e^{\frac{C_2 t}{2^{d/2}}} - 1), \quad (17)$$

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 7

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:

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

$$G(\theta) \leq G(\theta_0) + 2L\|\theta - \theta_0\|^2. \quad (18)$$

- Therefore,

$$\frac{d\mathbb{E}_f[\|\theta_t^f - \theta_0\|^2]}{dt} = 2\mathbb{E}_f[\langle \theta_t^f - \theta_0, -\nabla_{\theta} R^f(\theta_t^f) \rangle] \quad (19)$$

$$\leq \frac{1}{2^{\frac{d}{2}-1}} \sqrt{\mathbb{E}_f[\|\theta_t^f - \theta_0\|^2] \mathbb{E}_f[G(\theta_t^f)]} \quad (20)$$

$$\leq \frac{1}{2^{\frac{d}{2}-1}} \sqrt{\mathbb{E}_f[\|\theta_t^f - \theta_0\|^2] \mathbb{E}_f[G(\theta_0) + 2L\|\theta_t^f - \theta_0\|^2]}. \quad (21)$$

Hardness of learning with GD

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

$$\dot{\delta}_t \leq 2^{2-\frac{d}{2}}(\sqrt{2L}\delta_t + \sqrt{G(\theta_0)}). \quad (22)$$

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.

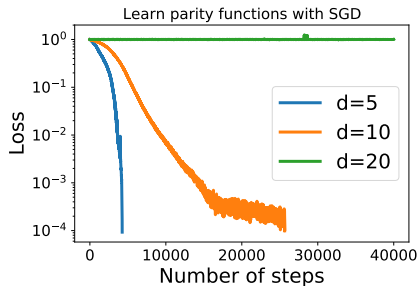
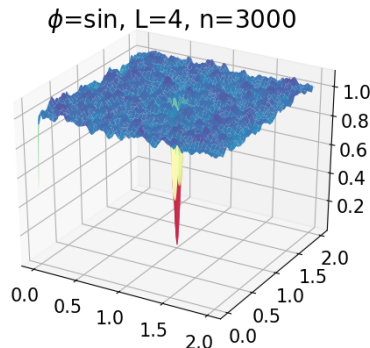
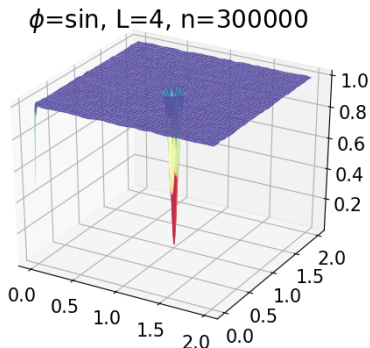


Figure 1: Learning parity functions with SGD and two-layer neural networks.

An illustration of the landscape of sine neuron

Consider a two-dimensional case. $R(\mathbf{w}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[(\sin(L2\pi\mathbf{w}^T \mathbf{x}) - \sin(L2\pi\mathbf{w}^T \mathbf{x}))^2]$, 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 variance (wrt the target function) of gradients is exponentially small. Hence, *gradient-based* algorithms are unlikely to succeed. This observation hold for any parametric model as long as it is satisfies 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 variance (wrt the target function) of gradients is exponentially small. Hence, *gradient-based* algorithms are unlikely to succeed. This observation holds for any parametric model as long as it satisfies some smooth condition.
- Typical examples include the parity function and the cosine neuron: $f(\mathbf{x}) = \cos(\mathbf{w}^T \mathbf{x})$. These functions can be represented as two-layer ReLU networks with the path norm bounded by Cd^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 variance (wrt the target function) of gradients is exponentially small. Hence, *gradient-based* algorithms are unlikely to succeed. This observation holds for any parametric model as long as it satisfies some smooth condition.
- Typical examples include the parity function and the cosine neuron: $f(\mathbf{x}) = \cos(\mathbf{w}^T \mathbf{x})$. These functions can be represented as two-layer ReLU networks with the path norm bounded by Cd^2 .
- These hardness results suggest that the Barron space is very likely too large to study the training of two-layer neural networks.