

Lecture 17: Computational Complexity I

Today, we introduced the next unit on computational complexity, specifically on proving lower bounds of problems. These lower bounds serve to tell us when continuing to find a more efficient algorithm for a problem is futile and when we should revisit and possibly alter the assumptions we made.

1 Finishing Lecture 16

1.1 Finishing Mean Field Limit

Recall that last time we had a MLP with scaling defined as

$$f_{\theta}(x) = \frac{1}{N} \sum_i a_i \sigma(\langle w_i, x \rangle). \quad (1)$$

Here N represents the “width” of the network (the number of neurons). We have previously observed that as $N \rightarrow \infty$, the neurons at time t of the gradient flow will converge to i.i.d. draws from the distribution ρ_t , which satisfies the differential equation

$$\partial_t \rho_t = \operatorname{div}(\rho_t \cdot \nabla \psi_{\rho_t}) \quad (2)$$

However, in many settings, this PDE is intractable to work with. Hence, in order to greatly simplify the problem, we often consider toy models, specifically ones with high symmetry. For example, one such simplified, model that makes use of symmetry is when the x 's are normally distributed and $y = \phi(\langle w^*, x \rangle)$. This is a type of single index model, where the data depends only on a fixed projection of the data (in the direction of w^*). Hence, in this model the data is secretly 1-dimensional. As in this situation, the data distribution doesn't vary with rotations so long as w^* is preserved, this greatly simplifies the PDE.

In these settings you can often numerically solve the PDE to obtain actual predictions. However, such numerical methods rely on the simplifying assumptions we made earlier. However, it is still difficult to get provable results from these algorithms. Currently, the only “end-to-end” non-asymptotic, global convergence

result for learning these types of single-index models with “standard” gradient descent is the following theorem, stated without proof:

Theorem 1. [MHD⁺23] *Projected gradient descent on a one-hidden-layer student network with quartic polynomial activations and polynomial width learns certain functions of the form $y = p(\langle w^*, x \rangle)$, where p is a degree-4 polynomial, over Gaussian inputs using $O(d^{3.1})$ samples.*

This theorem has an enormously complex proof that is beyond the scope of the class.

Now we will revisit the correlational statistical query model (CSQ). For single index-models (models of the form $y = \phi(\langle w^*, x \rangle)$), the complexity of CSQ algorithms are dictated primarily by what is known as the “informal exponent,” that is, the smallest s for which the s -th Hermite coefficient of ϕ is non-zero. If we are training only a single neuron, running online SGD will learn in an expected amount of time, $O(d^s)$. [AGJ21] There are results that generalize the idea of “leap complexity” to multi-index models (models of the form $y = \phi(\prod_W x)$) in time $O(d^{leap})$ [ABAM23].

In general, the algorithms discussed above are optimal for CSQ algorithms. However, there are (at least in principle) other more efficient non-CSQ algorithms, such as filtered PCA which achieve $O(d)$ sample complexity and fixed-polynomial runtime. The optimality of these CSQ algorithms, can be proved using computational lower-bounds (which we will discuss in a moment).

1.2 Recap of Supervised Learning

Still virtually all algorithms in the field of PAC learning, rely on low-degree polynomials. We began with approximating binary circuits by their low degree Fourier coefficients and ended with using the Mean Field Limit to learn low-degree components of the underlying functions generating the data.

However, there are two distinct concepts of low degree polynomials have discussed. The first is *What is the smallest degree for which there is a polynomial that approximates the ground truth?* This approach exploits Fourier/Hermite concentrations. Algorithms include low-degree algorithm, polynomial regression, and kernel methods.

The other concept of low-degree polynomials is *What is the smallest degree at which the ground truth has a nonzero polynomial component?* This approach exploits informa-

tion exponents/leaps. Algorithms included tensor methods, feature learning/GD beyond NTK.

2 Introduction to Computational Complexity

2.1 Guiding Examples

Throughout the entire course, we have consistently seen two examples show up:

- Learning a mixture of k Gaussians in \mathbb{R}^d
- Learning neural networks of size k over Gaussian *inputs* in \mathbb{R}^d

For both algorithms, it is an open question regarding the existence of a fully $\text{poly}(k, d)$ -time algorithm. This could either be because we just haven't found some mathematical solution yet, or it could be because no efficient algorithm exists as the problem is currently formulated. To rule out the first option, we make use of **lower bounds**.

2.2 Computational Hardness

We have dealt with several problems restricted by statistical lower bounds. For example, the Airy Disks from the first lecture will require exponentially many samples to differentiate below the diffraction limit. In this case, solving this problem without enough data is impossible even with an infinite amount of compute.

In this section, we will explore the case where there is enough signal in the dataset that a computationally *inefficient* algorithm (i.e. brute force) can solve the problem, but that no computationally efficient algorithm exists to solve the problem. There are a couple of different version of classical computational hardness such as NP hardness (SAT, graph coloring, etc.) and cartographic hardness (public-key cryptography, pseudorandom generators, etc.).

However, we will focus primarily on “average case hardness.” A classic example is the **planted clique** problem:

Setup: Given the adjacency matrix of a graph sampled either from an Erdos-Renyi graph $G(n, 1/2)$ (every edge independently included with probability $1/2$) or a Planted graph (graph is $G(n, 1/2)$ with a random clique of size N).

Task: Determine which case we are in with high probability over the randomness

of the instance.

The largest clique in $G(n, 1/2)$ is of size either $2 \log n$ or $2 \log n + 1$ with probability $1 - o_n(1)$. Thus, the problem is information-theoretically intractable when $N \in \Omega(\log n)$. We could brute force a solution to this algorithm (i.e. there is enough “signal”), but no known computationally efficient algorithm exists for $N = o(\sqrt{N})$. In the case of $N \in \Omega(\sqrt{N})$, there does exist an efficient algorithm using the top eigen vector of the adjacency matrix.[AKS98] It is hypothesized that no poly-time algorithm exists for this case.

Another problem is **learning parity with noise**, which is a noisy supervised learning task. For positive η , random $S \subseteq [d]$, and dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, defined as

$$x \sim \{\pm 1\}^d \quad y = \begin{cases} x_S & \text{w.p. } 1 - \eta \\ -x_S & \text{otherwise} \end{cases},$$

$N = O(d \log d)$ samples are information-theoretically sufficient, as we could (inefficiently) brute force the solution. When we have $\eta = 0$, this is just a linear system modulo 2 that we could solve in polynomial time with Gaussian elimination. It is hypothesized that no polynomial time algorithm exists even to distinguish the y 's from random labels.

2.3 Traditional Hardness Paradigm

The classic approach to proving hardness rests of reductions. Given some problem X that is provably hard, show that efficiently solving some other problem Y could be (efficiently) mapped to a solution of problem X . This approach is very effective for worst-case hardness but break down for average-case hardness. This requires us to look for new ways to prove average-case hardness that do not rely on reductions.

3 Lower Bounds in Restricted Models

3.1 Restricted Model of Computation

Here is an alternate approach to proving these lower-bounds:

- Formally define a restricted model of computation that captures all known algorithms for the problem in question.
- Prove a lower bound against algorithms in the restricted model

This approach works unconditionally (unlike NP hardness which is conditional on $P \neq NP$). There are a few different approaches to this proof:

- **Statistical query:** Only makes use of noisy estimates of population-level statistics of the data distribution
- **Lipschitz algorithms:** If the instance is perturbed, the algorithm output does not change much.
- **Low-degree algorithms:** Only uses low-degree polynomials evaluated on the data
- **Sum-of-Squares Algorithms:** There is a “Canonical” SoS relaxation of the problem, and we want to prove high degree SoS is necessary

These are much “easier” to show than reductions; however, they are often less convincing than reduction-based bounds.

3.2 Statistical Query

We will return to the correlational statistical query (CSQ) model of computation. We will define the statistical query (SQ) model of computation as an algorithm that only interacts with some dataset $\{(x_i, y_i)\}_{i=1}^N$ through some oracle that takes in

$$\psi : \mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}$$

and outputs

$$\mathbb{E}[\psi(x, y)] + \text{noise},$$

where $|\text{noise}| \leq \tau$ for some tolerance $\tau = \sqrt{1/N}$.

If the labels $y_i \in \{\pm 1\}$ and $x \sim q$ for some unknown distribution q , then $SQ = CSQ$ and we have

$$\mathbb{E}[\psi(x, y)] = \mathbb{E} \left[\frac{1+y}{2} \psi(x, 1) + \frac{1-y}{2} \psi(x, -1) \right] = \mathbb{E}[g(x)] + \mathbb{E}[y \cdot h(x)]. \quad (3)$$

In general, the statistical query model can capture essentially any known learning algorithm except for Gaussian elimination and others (it is still an open question which algorithms it cannot capture).

3.3 Proof of SQ Lower Bound for Noiseless Parity

Theorem 2. [Kea98] Any statistical query algorithm for learning parity (without noise) requires $2^{\Omega(d)}$ queries or tolerance $2^{-\Omega(d)}$.

Proof. Recall the setting from above for learning parity *with* noise. We have the data set $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ with

$$x \sim \{\pm 1\}^d \quad y = \begin{cases} x_S & \text{w.p. } 1 - \eta \\ -x_S & \text{otherwise} \end{cases}.$$

Now, however, for the setting *without* noise, we will set $\eta = 0$. Consider first any CSQ $\phi : \{\pm 1\}^d \rightarrow [-1, 1]$. Let

$$\phi_S := \mathbb{E}_x[x_S \phi(x)].$$

I claim first that for uniformly random subsets $S \subseteq [n]$,

$$\text{Var}_S[\phi_S] \geq 2^{-\Omega(n)}.$$

The proof of this lemma is as follows:

$$\begin{aligned} \text{Var}_S[\phi_S] &= \mathbb{E}_S[\phi_S^2] - \mathbb{E}_S[\phi_S]^2 \\ &= \mathbb{E}_S \mathbb{E}_{x, x'}[x_S x'_S \phi(x) \phi(x')] - \mathbb{E}_{S, S'} \mathbb{E}_{x, x'}[x_S x'_S \phi(x) \phi(x')] \\ &= \mathbb{E}_{x, x'}[\phi(x) \phi(x')] \underbrace{\mathbb{E}_{S, S'}[x_S x'_S - x_S x_{-S'}]}_{\alpha} \quad (*) \end{aligned}$$

I claim that the inside expectation α is equal to 0 if $x \neq x'$. For $z \neq \mathbf{1}$, $E_S[z_S] = 0$ so if $x \neq x'$, then $\mathbb{E}_S[x_S x'_S] = \mathbb{E}[z_S] = 0$ and $\mathbb{E}_{S, S'}[x_S x'_S] = \mathbb{E}_S[x_S] \mathbb{E}_{S'}[x'_{S'}] = 0$ because at most of $x, x' = \mathbf{1}$. This is a “pairwise independence” argument.

Returning, to the previous expression, we can use this fact to get

$$\begin{aligned} (*) &= \mathbb{E}_{x, x'}[\phi(x) \phi(x') \cdot \mathbb{1}[x = x']] \\ &= \frac{1}{2^n} \mathbb{E}_x[\phi(x)^2] \\ &\leq \frac{1}{2^n}. \end{aligned}$$

Thus, by Chebyshev's inequality, we can get that

$$\begin{aligned}\mathbb{P}_S[|\phi_S - \mathbb{E}_S[\phi_S]| \geq \tau] &\leq \frac{1}{\tau^2} \text{Var}_S(\phi_S) \\ &\leq \frac{1}{2^n \tau^2}\end{aligned}$$

Hence, in order to answer the CSQ ϕ , we can just output $\mathbb{E}_S[\phi_S]$. If the tolerance is equal to τ , this is accurate for $\frac{1}{2^n \tau^2}$ fraction of parity functions. That is, each CSQ only rules out *at most* $1/\tau^2$ many S 's. Thus, we need $2^n/\tau^2$ many queries. Letting $\tau := 2^{-n/2}$ completes the proof. \square

References

- [ABAM23] Emmanuel Abbe, Enric Boix-Adsera, and Theodor Misiakiewicz. Sgd learning on neural networks: leap complexity and saddle-to-saddle dynamics, 2023.
- [AGJ21] Gerard Ben Arous, Reza Gheissari, and Aukosh Jagannath. Online stochastic gradient descent on non-convex losses from high-dimensional inference. *Journal of Machine Learning Research*, 22(106):1–51, 2021.
- [AKS98] Noga Alon, Michael Krivelevich, and Benny Sudakov. Finding a large hidden clique in a random graph. *ETH Zurich*, 1998.
- [Kea98] Michael Kearns. Efficient noise-tolerant learning from statistical queries. *J. ACM*, 45(6):983–1006, nov 1998.
- [MHD⁺23] Arvind Mahankali, Jeff Z. Haochen, Kefan Dong, Margalit Glasgow, and Tengyu Ma. Beyond ntk with vanilla gradient descent: A mean-field analysis of neural networks with polynomial width, samples, and time, 2023.