
Does Preprocessing Help Training Over-parameterized Neural Networks?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep neural networks have achieved impressive performance in many areas. De-
2 signing a fast and provable method for training neural networks is a fundamental
3 question in machine learning.

4 The classical training method requires paying $\Omega(mnd)$ cost for both forward com-
5 putation and backward computation, where m is the width of the neural network,
6 and we are given n training points in d -dimensional space. In this paper, we pro-
7 pose two novel preprocessing ideas to bypass this $\Omega(mnd)$ barrier:

- 8 • First, by preprocessing the initial weights of the neural networks, we can
9 train the neural network in $\tilde{O}(m^{1-\Theta(1/d)}nd)$ cost per iteration.
- 10 • Second, by preprocessing the input data points, we can train neural network
11 in $\tilde{O}(m^{4/5}nd)$ cost per iteration.

12 From the technical perspective, our result is a sophisticated combination of tools
13 in different fields, greedy-type convergence analysis in optimization, sparsity ob-
14 servation in practical work, high-dimensional geometric search in data structure,
15 concentration and anti-concentration in probability. Our results also provide theo-
16 retical insights for a large number of previously established fast training methods.

17 1 Introduction

18 During the last decade, deep learning has achieved dominating performance over many ar-
19 eas, e.g., computer vision [LBBH98, KSH12, SLJ⁺15, HZRS16], natural language processing
20 [CWB⁺11, DCLT18], automatic driving system, game playing [SHM⁺16, SSS⁺17] and beyond.
21 The computational resource requirement for deep neural network training such computation has
22 growing very quickly. Designing a fast and provable training method for neural networks is there-
23 fore a fundamental and demanding challenge.

24 Almost all deep learning models are optimized by gradient descent (or its variants). The total training
25 time can be split in two components, the first one is the number of iterations and the second one is the
26 cost per spent per iteration. Nearly all the iterative algorithms for acceleration can be viewed as two
27 separate lines of research correspondingly, the first line is aiming for an algorithm that has as small
28 as possible number of iterations, the second line is focusing on designing as efficient as possible data
29 structures to improve the cost spent per iteration of the algorithm [Vai89, CLS19, LSZ19, JLSW20,
30 JKL⁺20]. In this paper, our major focus is on the second line.

31 There are a number of practical works trying to use a nearest neighbor search data structure to speed
32 up the per-step computation of the deep neural network training [CMF⁺20, LXJ⁺20, CLP⁺21,
33 DMZS21]. However, none of the previous work is able to give a provable guarantee. In this paper,
34 our goal is to develop training algorithms that provably reduce per step time complexity. Let us

consider the ReLU activation neural network and two-layer neural network¹. Let n denote the number of training data points. Let d denote the dimension of each data point. Let m denote the number of neurons. Suppose we consider the gradient descent algorithm, in each iteration we need to compute prediction for each point in the neural network. Each point $x_i \in \mathbb{R}^d$, requires to compute m inner product in d dimension. Thus, $\Omega(mnd)$ is a natural barrier for cost per iteration in training neural network (in both forward computation and backward computation).

A natural question to ask is

Is it possible to improve the cost per iteration of training neural network algorithm? E.g., is $o(mnd)$ possible?

We list our contributions as follows:

- We provide a new theoretical framework for speeding up neural network training by: 1) adopting the shifted neural tangent kernel; 2) showing that only a small fraction ($o(m)$) of neurons are activated for each input data in each training iteration; 3) identifying the sparsely activated neurons via geometric search; 4) proving that the algorithm can minimize the train loss to zero in a linear convergence rate.
- We provide two theoretical results 1) our first result (Theorem 6.1) builds a dynamic half-space report data structure for the weights of neural network to train neural networks in sublinear cost per iteration; 2) our second result (Theorem 6.2) builds a static half-space report data-structure for the input data points of the training data set for training neural network in sublinear time.

1.1 Related work

Acceleration via high-dimensional search data-structure. High-dimensional search data structures support efficiently finding points in some geometric query regions (e.g., half-spaces, simplices, etc). Currently, there are two main approaches: one is based on Locality Sensitive Hashing (LSH) [IM98], which aims to find the close-by points (i.e., small ℓ_2 distance [DIIM04, AR15, AIL⁺15, ARN17, Raz17, AIR18, BIW19, DIRW20] or large inner product [SL14, SL15b, SL15a]) of a query $q \in \mathbb{R}^d$ in a given set of points $S \subset \mathbb{R}^d$. This kind of algorithms runs very fast in practice, but most of them only supports approximate queries. Another approach is based on space partitioning data structures, for example, partition trees [Mat92a, Mat92b, AEM92, AC09, Cha12], k - d trees / range trees [CT17, TOG17, Cha19], Voronoi diagrams [ADBMS98, Cha00], which can exactly search the query regions. Recent works have successfully applied high-dimensional geometric data-structure to reduce the complexity of training deep learning models. SLIDE [CMF⁺20] accelerates the forward pass by retrieving neurons with maximum inner product via a LSH-based data-structure; Reforemer [KKL20] similarly adopts LSH to reduce the memory usage for processing long sequence; MON-GOOSE [CLP⁺21] accelerates the forward pass by retrieving neurons with maximum inner product via a learnable LSH-based data-structure [Cha02] and lazy update framework [CLS19]. Despite the great empirical success, there is no theoretical understanding of such acceleration.

The goal of our paper is to theoretically characterize the acceleration brought by the high-dimensional geometric data-structure. Specifically, our algorithm and analysis are built upon the HSR data structures [AEM92] which can find all the points that have large inner product and support efficient data update. Note that HSR comes with a stronger recovery guarantee than LSH, in the sense that HSR, whereas LSH is guaranteed to find some of those points.

Convergence via over-parameterization. Recently, there has been a tremendous progress in understanding the “small training loss” phenomenon in standard training [LL18, DZPS19, AZLS19a, AZLS19b, DLL⁺19, ADH⁺19a, ADH⁺19b, SY19, OS20, LSS⁺20, ZPD⁺20, HLSY21]. A convergence theory has been developed to show that, when randomly initialized, gradient descent and stochastic gradient descent converge to small training loss in polynomially many iterations when the network has polynomial width in terms of the number of training examples.

¹An alternative name of two-layer neural network is “one-hidden layer neural network”

2 Our techniques

- Empirical works combine high-dimensional search data structures (e.g., LSH) with neural network training, however they do not work theoretically due to the following reasons:
 - Without shifting, the number of activated (and therefore updated) number of neurons is $\Theta(m)$. There is no hope to theoretically prove $o(m)$ complexity (See **Issue 1**).
 - Approximate high-dimensional search data structure might miss some important neurons, which can potentially prevent the training from converge (see **Issue 2**).
- Our solutions are:
 - We propose a shifted ReLU activation that is guaranteed to have $o(m)$ number of activated neurons. Along with the shifted ReLU, we also propose a shifted NTK to rigorously provide convergence guarantee (see **Solution 1**).
 - We adopt an exact high-dimensional search data structure which better couples with the shifted NTK. It takes $o(m)$ time to identify the activated neurons and fits well with the convergence analysis as it avoids missing important neurons (see **Solution 2**).

Issue 1 To speed up the training process, we need the neural network to be “sparse”, that is, for each training data $x \in \mathbb{R}^d$, the number of activated neurons is small. Then, in the forward computation, we can just evaluate a small subset of neurons. However, in the previous NTK analysis (e.g., [DZPS19]), the activation function is $\sigma(x) = \max\{\langle w_r, x \rangle, 0\}$, and the weights vectors w_r are initially sampled from a standard d -dimensional Gaussian distribution. Then, by the symmetry of Gaussian distribution, we know that for every input data x , there will be about half of the neurons being activated, which means that we can only obtain a constant factor speedup.

Solution 1 The problem actually comes from the activation function. In practice, people use a shifted ReLU function $\sigma_b(x) = \max\{\langle w_r, x \rangle, b_r\}$ to train neural networks. The main observation of our work is that *threshold implies sparsity*. We consider the setting where all neurons have a unified threshold parameter b . Then, by the concentration of Gaussian distribution, there will be $O(\exp(-b^2) \cdot m)$ activated neurons after the initialization.

The next step is to show that the number of activated neurons will not blow up too much in the following training iterations. In [DZPS19, SY19], they showed that the weights vectors are changing slowly during the training process. In our work, we open the black box of their proof and show a similar phenomenon for the shifted ReLU function. More specifically, a key component is to prove that for each training data, a large fraction of neurons will not change their status (from non-activated to activated and vice versa) in the next iteration with high probability. To achieve this, they showed that this is equivalent to the event that a standard Gaussian random variable in a small centered interval $[-R, R]$, and applied the anti-concentration inequality to upper-bound the probability. In our setting, we need to upper-bound the probability of $z \sim \mathcal{N}(0, 1)$ in a shifted interval $[b - R, b + R]$. On the one hand, we can still apply the anti-concentration inequality by showing that the probability is at most $\Pr[z \in [-R, R]]$. On the other hand, this probability is also upper-bounded by $\Pr[z > b - R]$, and for small R , we can apply the concentration inequality for a more accurate estimation. In the end, by some finer analysis of the probability, we can show that with high probability, the number of activated neurons in each iteration is also $O(\exp(-b^2) \cdot m)$ for each training data. If we take $b = \Theta(\sqrt{\log m})$, we only need to deal with truly sublinear in m of activated neurons in the forward evaluation.

Issue 2: How to find the small subset of activated neurons? A linear scan of the neurons will lead to a time complexity linear in m , which we hope to avoid. Randomly sampling or using LSH for searching can potentially miss important neurons which are important for a rigorous convergence analysis.

Solution 2 Given the shifted ReLU function $\sigma_b(\langle w_r, x \rangle) = \max\{\langle w_r, x \rangle - b, 0\}$, the active neurons are those with weights w_r lying in the half space of $\langle w_r, x \rangle - b > 0$. Finding such neurons is equivalent to a computational geometry problem: given m points in \mathbb{R}^d , in each query and a half space \mathcal{H} , the goal is to output the points contained in \mathcal{H} . Here we use the Half-Space Reporting (HSR) data structure proposed by [AEM92]: after proper initialization, the HSR data structure can return all points lying in the queried half space with complexity as low as $O(\log(n) + k)$, where k

is the number of such points. Note that the HSR data structure well couples with the shifted ReLU, as the number of activated neurons k is truly sublinear in m as per the setting of $b = \Theta(\sqrt{\log m})$.

3 Preliminaries

Notations For an integer n , we use $[n]$ to denote the set $\{1, 2, \dots, n\}$. For a vector x , we use $\|x\|_2$ to denote the entry-wise ℓ_2 norm of a vector. For a matrix A , we use $\|A\|_F = (\sum_{i,j} A_{i,j}^2)^{1/2}$ to denote its Frobenius norm and use $\|A\|$ to denote the operator/spectral norm of A . We use $x^\top y$ to denote the inner product between vectors x and y . We use I_d to denote d -dimensional identity matrix. We use $\mathcal{N}(0, 1)$ to denote Gaussian distribution with mean 0 and variance 1.

3.1 Problem Formulation

In this section, we introduce the neural network model we study in this paper. We consider a two-layer ReLU activated neural network with m neurons in the hidden layer²:

$$f(W, x, a) := \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma_b(\langle w_r, x \rangle),$$

where $x \in \mathbb{R}^d$ is the input, $w_1, \dots, w_m \in \mathbb{R}^d$ are weight vectors in the first layer, $a_1, \dots, a_m \in \mathbb{R}$ are weights in the second layer. The ReLU function $\sigma_b(x) := \max\{x - b, 0\}$, where b is the threshold parameter. For simplicity, we only optimize $W \in \mathbb{R}^{d \times m}$ but not $a \in \mathbb{R}^m$.³ Specifically, we use the following initialization throughout the paper

$$w_r(0) \sim \mathcal{N}(0, I_d); \quad a_r \sim \mathcal{U}(\{-1, 1\}), \quad \forall r \in [m]. \quad (1)$$

Recall that the ReLU function $\sigma_b(x) = \max\{x - b, 0\}$. Therefore for $r \in [m]$, we have

$$\frac{\partial f(W, x, a)}{\partial w_r} = \frac{1}{\sqrt{m}} a_r x \mathbf{1}_{w_r^\top x \geq b}. \quad (2)$$

We define objective function L as $L(W) = \frac{1}{2} \sum_{i=1}^n (y_i - f(W, x_i, a))^2$.

Then, we can apply the gradient descent to optimize the weight matrix W in the following standard way,

$$W(k+1) = W(k) - \eta \frac{\partial L(W(k))}{\partial W(k)}. \quad (3)$$

We can compute the gradient of L in terms of $w_r \in \mathbb{R}^d$

$$\frac{\partial L(W)}{\partial w_r} = \frac{1}{\sqrt{m}} \sum_{i=1}^n (f(W, x_i, a) - y_i) a_r x_i \mathbf{1}_{\langle w_r, x_i \rangle \geq 0}. \quad (4)$$

At time t , let $u(t) = (u_1(t), \dots, u_n(t)) \in \mathbb{R}^n$ be the prediction vector where each $u_i(t)$ is defined as

$$u_i(t) = f(W(t), a, x_i). \quad (5)$$

3.2 Data Structure for Half-Space Reporting

The half-space range reporting problem is an important problem in computational geometry, which is formally defined as follows:

Definition 3.1 (Half-space range reporting). *Given a set S of n points in \mathbb{R}^d . There are two operations:*

²This is a very standard formulation in the literature, e.g., see [DZPS19, SY19].

³We remark, in some previous work, they do choose shift, but their shift is a random shift. In our application, it is important, we fixed the same b for all neurons and never trained.

- 161 • QUERY(H): given a half-space $H \subset \mathbb{R}^d$, output all of the points in S that contain in H ,
162 i.e., $S \cap H$.
- 163 • UPDATE: add or delete a point in S .
- 164 – INSERT(q): insert q into S
- 165 – DELETE(q): delete q from S

166 Let $\mathcal{T}_{\text{init}}$ denote the pre-processing time to build the data structure, $\mathcal{T}_{\text{query}}$ denote the time per query
167 and $\mathcal{T}_{\text{update}}$ time per update.

168 We use the data-structure proposed in [AEM92] to solve the half-space range reporting problem,
169 which admits the interface summarized in Algorithm 1. Intuitively, the data-structure recursively
170 partition the set S and organizes the points in a tree data-structure. Then for a given query (a, b) , all
171 k points of S with $\text{sgn}(\langle a, x \rangle - b) \geq 0$ are reported quickly. Note that the query (a, b) here defines
172 the half-space H in Definition 3.1.

Algorithm 1 Half Space Report Data Structure

```

1: data structure HALFSpacERePORT
2:   procedures:
3:     INIT( $S, n, d$ )           ▷ Initialize the data structure with a set  $S$  of  $n$  points in  $\mathbb{R}^d$ 
4:     QUERY( $a, b$ )             ▷  $a, b \in \mathbb{R}^d$ . Output the set  $\{x \in S : \text{sgn}(\langle a, x \rangle - b) \geq 0\}$ 
5:     ADD( $x$ )                  ▷ Add a point  $x \in \mathbb{R}^d$  to  $S$ 
6:     DELETE( $x$ )               ▷ Delete the point  $x \in \mathbb{R}^d$  from  $S$ 
7:   end data structure

```

173 Adapted from [AEM92], the algorithm comes with the following complexity:

174 **Corollary 3.2** ([AEM92]). Given a set of n points in \mathbb{R}^d , the half-space reporting problem can be
175 solved with the following performances:

- 176 • Part 1.⁴ $\mathcal{T}_{\text{query}}(n, d, k) = O_d(n^{1-1/\lfloor d/2 \rfloor} + k)$, amortized $\mathcal{T}_{\text{update}} = O_d(\log^2(n))$.
- 177 • Part 2.⁵ $\mathcal{T}_{\text{query}}(n, d, k) = O_d(\log(n) + k)$, amortized $\mathcal{T}_{\text{update}} = O_d(n^{\lfloor d/2 \rfloor - 1})$.

178 3.3 Sparsity-based Characterizations

179 In this section, we consider the ReLU function with a nonzero threshold: $\sigma_b(x) = \max\{0, x - b\}$,
180 which is commonly seen in practise, and also has been considered in theoretical work [ZPD⁺20].

181 We first define the set of neurons that are firing at time t .

182 **Definition 3.3** (fire set). For each $i \in [n]$, for each $t \in \{0, 1, \dots, T\}$, let $\mathcal{S}_{i, \text{fire}}(t) \subset [m]$ denote
183 the set of neurons that are “fire” at time t , i.e.,

$$\mathcal{S}_{i, \text{fire}}(t) := \{r \in [m] : \langle w_r(t), x_i \rangle > b\}.$$

184 We define $k_{i,t} := |\mathcal{S}_{i, \text{fire}}(t)|$, $\forall t \in \{0, 1, \dots, T\}$.

185 The following lemma shows that σ_b gives the desired sparsity.

186 **Lemma 3.4** (Sparsity after initialization). Let $b > 0$ be a tunable parameter. If we use the σ_b as
187 the activation function, then after the initialization, with probability at least $1 - n \cdot \exp(-\Omega(m \cdot$
188 $\exp(-b^2/2)))$, it holds that for each input data x_i , the number of activated neurons $k_{i,0} \leq O(m \cdot$
189 $\exp(-b^2/2))$, where m is the total number of neurons.

190 *Proof.* By the concentration of Gaussian distribution, the initial fire probability of a single neuron is

$$\Pr[\sigma_b(\langle w_r(0), x_i \rangle) > 0] = \Pr_{z \sim \mathcal{N}(0,1)}[z > b] \leq \exp(-b^2/2).$$

⁴Used in Theorem 6.1

⁵Used in Theorem 6.2

Algorithm 2 Training Neural Network via building a data structure of weights of the neural network

```

1: procedure TRAININGWITHPREPROCESSWEIGHTS( $\{(x_i, y_i)\}_{i \in [n]}, n, m, d$ ) ▷ Theorem 6.1
2:   Initialize  $w_r, a_r$  for  $r \in [m]$  and  $b$  according to Equation (1) and Remark 3.5
3:   HALFSACEREPOR HSR.INIT( $\{w_r(0)\}_{r \in [m]}, m, d$ ) ▷ Algorithm 1
4:   for  $t = 1 \rightarrow T$  do
5:      $S_{i, \text{fire}} \leftarrow \text{HSR.QUERY}(x_i, b)$  for  $i \in [n]$ 
6:     Forward pass for  $x_i$  only on neurons in  $S_{i, \text{fire}}$  for  $i \in [n]$ 
7:     Calculate gradient for  $x_i$  only on neurons in  $S_{i, \text{fire}}$  for  $i \in [n]$ 
8:     Gradient update for the neurons in  $\cup_{i \in [n]} S_{i, \text{fire}}$ 
9:     HSR.DELETE( $w_r(t)$ ) for  $r \in \cup_{i \in [n]} S_{i, \text{fire}}$ 
10:    HSR.ADD( $w_r(t+1)$ ) for  $r \in \cup_{i \in [n]} S_{i, \text{fire}}$ 
11:  end for
12:  return Trained weights  $w_r(T+1)$  for  $r \in [m]$ 
13: end procedure

```

191 Hence, for the indicator variable $\mathbf{1}_{r \in S_{i, \text{fire}}(0)}$, we have

$$\mathbb{E}[\mathbf{1}_{r \in S_{i, \text{fire}}(0)}] \leq \exp(-b^2/2).$$

192 By Bernstein inequality (Lemma B.3) we have for all $t > 0$,

$$\Pr[|S_{i, \text{fire}}(0)| > k_0 + t] \leq \exp\left(-\frac{t^2/2}{k_0 + t/3}\right), \quad (6)$$

193 where $k_0 := m \cdot \exp(-b^2/2)$. If we choose $t = k_0$, then we have:

$$\Pr[|S_{i, \text{fire}}(0)| > 2k_0] \leq \exp(-3k_0/8)$$

194 Then, by union bound over all $i \in [n]$, we have that with high probability

$$1 - n \cdot \exp(-\Omega(m \cdot \exp(-b^2/2))),$$

195 the number of initial fire neurons for the sample x_i is bounded by $k_{i,0} \leq 2m \cdot \exp(-b^2/2)$. □

196 The following remark gives an example of setting the threshold b , and will be useful for showing the
 197 sublinear complexity in the next section.

198 **Remark 3.5.** If we choose $b = \sqrt{0.4 \log m}$ then $k_0 = m^{4/5}$. For $t = m^{4/5}$, Eq. (6) implies that

$$\Pr[|S_{i, \text{fire}}(0)| > 2m^{4/5}] \leq \exp\left(-\min\{mR, O(m^{4/5})\}\right).$$

199 4 Training Neural Network with Half-Space Reporting Data Structure

200 In this section, we present two sublinear time algorithms for training over-parameterized neural net-
 201 works. The first algorithm (Section 4.1) relies on building a high-dimensional search data-structure
 202 for the weights of neural network. The second algorithm (Section 4.2) is based on building a data
 203 structure for the input data points of the training set. Both of the algorithms use the HSR to quickly
 204 identify the fired neurons to avoid unnecessary calculation. The time complexity and the sketch of
 205 the proof are provided after each of the algorithms.

206 4.1 Weights Preprocessing

207 We first introduce the algorithm that preprocesses the weights w_r for $r \in [m]$, which is commonly
 208 used in practice [CLP⁺21, CMF⁺20, KKL20]. Recall the two-layer ReLU neural network is defined
 209 as $f(W, x, a) := \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \sigma_b(\langle w_r, x \rangle)$. By constructing a HSR data-structure for w_r s, we can
 210 quickly find the set of active neurons $S_{i, \text{fire}}$ for each of the training sample x_i . See pseudo-code in
 211 Algorithm 2.

212 In the remaining part of this section, we focus on the time complexity analysis of Algorithm 2. The
 213 convergence proof will be given in Section 5.

Algorithm 3 Training Neural Network via building a data-structure of the training data

```

1: procedure TRAININGWITHPROCESSDATA( $\{(x_i, y_i)\}_{i \in [n]}, n, m, d$ ) ▷ Theorem 6.2
2:   Initialize  $w_r, a_r$  for  $r \in [m]$  and  $b$  according to Equation (1) and Remark 3.5
3:   HALFSACEREPOR HSR.INIT( $\{x_i\}_{i \in [n]}, n, d$ ) ▷ Algorithm 1
4:    $\tilde{S}_{r, \text{fire}} \leftarrow \text{HSR.QUERY}(w_r(0), b)$  for  $r \in [m]$  ▷  $\tilde{S}_{r, \text{fire}}$  are samples which neuron  $r$  fires for
5:    $S_{i, \text{fire}} \leftarrow \{r \mid i \in \tilde{S}_{r, \text{fire}}\}$  ▷  $S_{i, \text{fire}}$  is the set of neurons, which fire for  $x_i$ 
6:   for  $t = 1 \rightarrow T$  do
7:     Forward pass for  $x_i$  only on neurons in  $S_{i, \text{fire}}$  for  $i \in [n]$ 
8:     Calculate gradient for  $x_i$  only on neurons in  $S_{i, \text{fire}}$  for  $i \in [n]$ 
9:     Gradient update for the neurons in  $\cup_{i \in [n]} S_{i, \text{fire}}$ 
10:    for  $r \in \cup_{i \in [n]} S_{i, \text{fire}}$  do
11:       $S_{i, \text{fire}}.\text{DEL}(r)$  for  $i \in \tilde{S}_{r, \text{fire}}$ 
12:       $\tilde{S}_{r, \text{fire}} \leftarrow \text{HSR.QUERY}(w_r(t+1), b)$ 
13:       $S_{i, \text{fire}}.\text{ADD}(r)$  for  $i \in \tilde{S}_{r, \text{fire}}$ 
14:    end for
15:  end for
16:  return Trained weights  $w_r(T+1)$  for  $r \in [m]$ 
17: end procedure

```

214 **Lemma 4.1** (Running time part of Theorem 6.1). *Given n data points in d -dimensional space.*
 215 *Running gradient descent algorithm (Algorithm 2) on a two-layer ReLU (over-parameterized) neu-*
 216 *ral network with $b = \sqrt{0.4 \log m}$, the expected per-iteration running time of the gradient descent*
 217 *algorithm is $\tilde{O}(m^{1-\Theta(1/d)}nd)$.*

218 *Proof.* The per-step time complexity is

$$\sum_{i=1}^n \mathcal{T}_{\text{QUERY}}(m, d, k_{i,t}) + (\mathcal{T}_{\text{DELETE}} + \mathcal{T}_{\text{INSERT}}) \cdot |\cup_{i \in [n]} S_{i, \text{fire}}(t)| + d \sum_{i \in [n]} k_{i,t}$$

219 The first term $\sum_{i=1}^n \mathcal{T}_{\text{QUERY}}(m, d, k_{i,t})$ corresponds to the running time of querying the active neuron
 220 set $S_{i, \text{fire}}(t)$ for all training samples $i \in [n]$. With the first result in Corollary B.7, the complexity is
 221 bounded by $\tilde{O}(m^{1-\Theta(1/d)}nd)$.

222 The second term $(\mathcal{T}_{\text{DELETE}} + \mathcal{T}_{\text{INSERT}}) \cdot |\cup_{i \in [n]} S_{i, \text{fire}}(t)|$ corresponds to updating w_r in the high-
 223 dimensional search data-structure (Lines 9 and 10). Again with the first result in Corollary B.7, we
 224 have $\mathcal{T}_{\text{DELETE}} + \mathcal{T}_{\text{INSERT}} = O(\log^2 m)$. Combining with the fact that $|\cup_{i \in [n]} S_{i, \text{fire}}(t)| \leq |\cup_{i \in [n]} S_{i, \text{fire}}(0)| \leq O(nm^{4/5})$, the second term is bounded by $O(nm^{4/5} \log^2 m)$.

226 The third term is the time complexity of gradient calculation restricted to the set $S_{i, \text{fire}}(t)$. With the
 227 bound on $\sum_{i \in [n]} k_{i,t}$ (Lemma C.9), we have $d \sum_{i \in [n]} k_{i,t} \leq O(m^{4/5}nd)$.

228 Putting them together completes the proof. □

229 4.2 Data Preprocessing

230 While the weights preprocessing algorithm is inspired by the common practise, the dual relationship
 231 between the input x_i and model weights w_r inspires us to preprocess the dataset before training (i.e.,
 232 building HSR data-structure for x_i). This largely improves the per-iteration complexity and avoids
 233 the frequent updates of the data structure since the training data is fixed. More importantly, once
 234 the training dataset is preprocessed, it can be reused for different models or tasks, thus one does not
 235 need to perform the expensive preprocessing for each training.

236 The corresponding pseudocode is presented in Algorithm 3. With x_i preprocessed, we can query
 237 HSR with weights w_r and the result $\tilde{S}_{r, \text{fire}}$ is the set of training samples x_i for which w_r fires for.
 238 Given $\tilde{S}_{r, \text{fire}}$ for $r \in [m]$, we can easily reconstruct the set $S_{i, \text{fire}}$, which is the set of neurons fired
 239 for sample x_i . The forward and backward pass can then proceed similar to Algorithm 2.

At the end of each iteration, we will update $\tilde{S}_{r,\text{fire}}$ based on the new w_r estimation and update $S_{i,\text{fire}}$ accordingly. For Algorithm 3, the HSR data-structure is static for the entire training process. This is the main difference from Algorithm 2, where the HSR needs to be updated every time step to account for the changing weights w_r .

We defer the convergence analysis to Section 5 and focus on the time complexity analysis of Algorithm 2 in the rest of this section. We consider d being a constant for the rest of this subsection.

Lemma 4.2 (Running time part of Theorem 6.2). *Given n data points in d -dimensional space. Running gradient descent algorithm (Algorithm 2) on a two-layer ReLU (over-parameterized) neural network with $b = \sqrt{0.4 \log m}$, the expected per-iteration running time of initializing $\tilde{S}_{r,\text{fire}}, S_{i,\text{fire}}$ for $r \in [m], i \in [n]$ is $O(m \log n + m^{4/5}n)$. The cost per iteration of the training algorithm is $O(m^{4/5}n \log n)$.*

Proof. We analyze the initialization and training parts separately.

Initialization In Lines 4 and 5, the sets $\tilde{S}_{r,\text{fire}}, S_{i,\text{fire}}$ for $r \in [m], i \in [n]$ are initialized. For each $r \in [m]$, we need to query the data structure the set of data points x 's such that $\sigma_b(w_r(0)^\top x) > 0$. Hence, the running time of this step is

$$\begin{aligned} \sum_{r=1}^m \mathcal{T}_{\text{query}}(n, d, \tilde{k}_{r,0}) &= O(m \log n + \sum_{r=1}^m \tilde{k}_{r,0}) \\ &= O(m \log n + \sum_{i=1}^n k_{i,0}) \\ &= O(m \log n + m^{4/5}n). \end{aligned}$$

where the second step follows from $\sum_{r=1}^m \tilde{k}_{r,0} = \sum_{i=1}^n k_{i,0}$.

Training Consider training the neural network for T steps. For each step, first notice that the forward and backward computation parts (Line 7 - 9) are the same as previous algorithm. The time complexity is $O(m^{4/5}n \log n)$.

We next show that maintaining $\tilde{S}_{r,\text{fire}}, r \in [m]$ and $S_{i,\text{fire}}, i \in [n]$ (Line 10-14) takes $O(m^{4/5}n \log n)$ time. For each fired neuron $r \in [m]$, we first remove the indices of data in the sets $S_{i,\text{fire}}$, which takes time

$$O(1) \cdot \sum_{r \in \cup_{i \in [n]} S_{i,\text{fire}}} \tilde{k}_{r,t} = O(1) \cdot \sum_{r=1}^m \tilde{k}_{r,t} = O(m^{4/5}n).$$

Then, we find the new set of x 's such that $\sigma_b(\langle w_r(t+1), x \rangle) > 0$ by querying the half-space reporting data structure. The total running time for all fired neurons is

$$\sum_{r \in \cup_{i \in [n]} S_{i,\text{fire}}} \mathcal{T}_{\text{query}}(n, d, \tilde{k}_{r,t+1}) \lesssim m^{4/5}n \log n + \sum_{r \in \cup_{i \in [n]} S_{i,\text{fire}}} \tilde{k}_{r,t+1} = O(m^{4/5}n \log n)$$

Then, we update the index sets $S_{i,\text{fire}}$ in time $O(m^{4/5}n)$. Therefore, each training step takes $O(m^{4/5}n \log n)$ time, which completes the proof. \square

5 Convergence Analysis of Our Algorithm

In this section, we state the convergence result of our training neural network algorithms (Lemma 5.2). An important component in our proof is to show that the smallest eigenvalue of the continuous Hessian matrix $\lambda_{\min}(H^{\text{cts}})$ cannot be too small. It turns out to be an anti-concentration problem of Gaussian random matrix. In [OS20], they gave a lower bound on $\lambda_{\min}(H^{\text{cts}})$ for ReLU function with $b = 0$, assuming the input data are separable. One of our major technical contribution is generalizing it to arbitrary $b \geq 0$.

273 **Proposition 5.1** (Informal version of Theorem F.1). *Given n input data $\{x_1, x_2, \dots, x_n\} \subseteq \mathbb{R}^d$*
 274 *such that $\forall i \in [n], \|x_i\|_2 = 1$ and $\delta := \min_{i \neq j} \{\|x_i - x_j\|_2, \|x_i + x_j\|_2\}$. For any $b \geq 0$, we define*
 275 *$H^{\text{cts}} \in \mathbb{R}^{n \times n}$ as follows $H_{i,j}^{\text{cts}} := \mathbb{E}_{w \sim \mathcal{N}(0, I_d)} [\langle x_i, x_j \rangle \mathbf{1}_{\langle w, x_i \rangle \geq b, \langle w, x_j \rangle \geq b}]$, $\forall i \in [n], j \in [n]$. Then*

$$\lambda_{\min}(H^{\text{cts}}) \geq 0.01e^{-b^2/2}\delta/n^2.$$

276 With proposition 5.1, we are ready to show the convergence rate of over-parameterized neural net-
 277 work with shifted ReLU function.

278 **Lemma 5.2** (Convergence part of Theorem 6.1 and Theorem 6.2). *Suppose input data-points are*
 279 *δ -separable, i.e., $\delta := \min_{i \neq j} \{\|x_i - x_j\|_2, \|x_i + x_j\|_2\}$. Let $m = \text{poly}(n, 1/\delta, \log(n/\rho))$, we*
 280 *i.i.d. initialize $w_r \in \mathcal{N}(0, I_d)$, a_r sampled from $\{-1, +1\}$ uniformly at random for $r \in [m]$,*
 281 *$b = \Theta(\sqrt{\log m})$, and we set the step size $\eta = O(\lambda/n^2)$ then with probability at least $1 - \rho$ over the*
 282 *random initialization we have for $k = 0, 1, 2, \dots, T$ ⁶,*

$$\|u(k) - y\|_2^2 \leq (1 - \eta\lambda/2)^k \cdot \|u(0) - y\|_2^2. \quad (7)$$

283 This result shows that despite the shifted ReLU and sparsely activated neurons, we can still retain
 284 the linear convergence. Combined with the results on per-step complexity in the previous section, it
 285 gives our main theoretical results of training deep learning models with sublinear time complexity
 286 (Theorem 6.1 and Theorem 6.2).

287 6 Main Theorems

288 In this section, we state the main theorems of our work, showing the sublinear running time and
 289 linear convergence rate of our two algorithms. The first algorithm is relying on building a high-
 290 dimensional geometric search data-structure for the weights of neural network.

291 **Theorem 6.1** (Main result I, informal of Theorem E.2). *Given n data points in d -dimensional space.*
 292 *We preprocess the initialization weights of the neural network. Running gradient descent algorithm*
 293 *(Algorithm 2) on a two-layer ReLU (over-parameterized) neural network with m neurons in the*
 294 *hidden layers is able to minimize the training loss to zero and the expected running time of gradient*
 295 *descent algorithm (per iteration)*

$$\tilde{O}(m^{1-\Theta(1/d)}nd).$$

296 The second algorithm is based on building a data structure for the input data points of the training
 297 set. Our second algorithm can further reduce the cost per iteration from $m^{1-1/d}$ to truly sublinear
 298 in m , e.g. $m^{4/5}$.

299 **Theorem 6.2** (Main result II, informal of Theorem E.2). *Given n data points in d -dimensional*
 300 *space. We preprocess all the data points. Running gradient descent algorithm (Algorithm 3) on a*
 301 *two-layer ReLU (over-parameterized) neural network with m neurons in the hidden layers is able*
 302 *to minimize the training loss to zero, and the expected running time of gradient descent algorithm*
 303 *(per iteration)*

$$\tilde{O}(m^{4/5}nd).$$

304 7 Discussion and Limitations

305 In this paper, we propose two sublinear algorithms to train neural networks. By preprocessing
 306 the weights of the neuron networks or preprocessing the training data, we rigorously prove that
 307 it is possible to train a neuron network with sublinear complexity, which overcomes the $\Omega(mnd)$
 308 barrier in classical training methods. Our results also offer theoretical insights for many previously
 309 established fast training methods.

310 One limitation of our work is that the current analysis framework does not provide convergence
 311 guarantee for combining LSH with gradient descent, which is commonly seen in many empirical
 312 works. Our proof breaks as LSH might miss important neurons which potentially ruins the conver-
 313 gence analysis. Instead, we refer to the HSR data-structure, which provides a stronger theoretical
 314 guarantee of successfully finding all fired neurons.

⁶Eventually, we choose $T = \lambda^{-2}n^2 \log(n/\epsilon)$ where ϵ is the final accuracy.

315 References

- 316 [AC09] Peyman Afshani and Timothy M Chan. Optimal halfspace range reporting in three di-
317 mensions. In *Proceedings of the twentieth annual ACM-SIAM symposium on Discrete*
318 *algorithms*, pages 180–186. SIAM, 2009.
- 319 [ADBMS98] Pankaj K Agarwal, Mark De Berg, Jiri Matousek, and Otfried Schwarzkopf. Con-
320 structing levels in arrangements and higher order voronoi diagrams. *SIAM journal on*
321 *computing*, 27(3):654–667, 1998.
- 322 [ADH⁺19a] Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained
323 analysis of optimization and generalization for overparameterized two-layer neural
324 networks. In *International Conference on Machine Learning*, pages 322–332, 2019.
- 325 [ADH⁺19b] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Ruslan Salakhutdinov, and Ruosong
326 Wang. On exact computation with an infinitely wide neural net. In *NeurIPS*. <https://arxiv.org/pdf/1904.11955.pdf>, 2019.
327
- 328 [AEM92] Pankaj K Agarwal, David Eppstein, and Jiri Matousek. Dynamic half-space report-
329 ing, geometric optimization, and minimum spanning trees. In *Annual Symposium on*
330 *Foundations of Computer Science (FOCS)*, volume 33, pages 80–80, 1992.
- 331 [AIL⁺15] Alexandr Andoni, Piotr Indyk, TMM Laarhoven, Ilya Razenshteyn, and Ludwig
332 Schmidt. Practical and optimal lsh for angular distance. In *Advances in Neural Infor-*
333 *mation Processing Systems (NIPS)*, pages 1225–1233. Curran Associates, 2015.
- 334 [AIR18] Alexandr Andoni, Piotr Indyk, and Ilya Razenshteyn. Approximate nearest neighbor
335 search in high dimensions. *arXiv preprint arXiv:1806.09823*, 7, 2018.
- 336 [AR15] Alexandr Andoni and Ilya Razenshteyn. Optimal data-dependent hashing for approx-
337 imate near neighbors. In *Proceedings of the forty-seventh annual ACM symposium on*
338 *Theory of computing (STOC)*, pages 793–801, 2015.
- 339 [ARN17] Alexandr Andoni, Ilya Razenshteyn, and Negev Shekel Nosatzki. Lsh forest: Practical
340 algorithms made theoretical. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM*
341 *Symposium on Discrete Algorithms (SODA)*, pages 67–78. SIAM, 2017.
- 342 [AZLS19a] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learn-
343 ing via over-parameterization. In *ICML*, 2019.
- 344 [AZLS19b] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. On the convergence rate of training
345 recurrent neural networks. In *NeurIPS*, 2019.
- 346 [Ber24] Sergei Bernstein. On a modification of chebyshev’s inequality and of the error formula
347 of laplace. *Ann. Sci. Inst. Sav. Ukraine, Sect. Math*, 1(4):38–49, 1924.
- 348 [BIW19] Arturs Backurs, Piotr Indyk, and Tal Wagner. Space and time efficient kernel density
349 estimation in high dimensions. 2019.
- 350 [Cha00] Timothy M Chan. Random sampling, halfspace range reporting, and construction of
351 ($\leq k$)-levels in three dimensions. *SIAM Journal on Computing*, 30(2):561–575, 2000.
- 352 [Cha02] Moses S Charikar. Similarity estimation techniques from rounding algorithms. In *Pro-*
353 *ceedings of the thirty-fourth annual ACM symposium on Theory of computing (STOC)*,
354 pages 380–388, 2002.
- 355 [Cha12] Timothy M Chan. Optimal partition trees. *Discrete & Computational Geometry*,
356 47(4):661–690, 2012.
- 357 [Cha19] Timothy M Chan. Orthogonal range searching in moderate dimensions: kd trees and
358 range trees strike back. *Discrete & Computational Geometry*, 61(4):899–922, 2019.
- 359 [Che52] Herman Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based
360 on the sum of observations. *The Annals of Mathematical Statistics*, pages 493–507,
361 1952.

362 [CLP⁺21] Beidi Chen, Zichang Liu, Binghui Peng, Zhaozhuo Xu, Jonathan Lingjie Li, Tri Dao,
363 Zhao Song, Anshumali Shrivastava, and Christopher Re. Mongoose: A learnable lsh
364 framework for efficient neural network training. In *ICLR oral*, 2021.

365 [CLS19] Michael B Cohen, Yin Tat Lee, and Zhao Song. Solving linear programs in the current
366 matrix multiplication time. In *STOC*, 2019.

367 [CMF⁺20] Beidi Chen, Tharun Medini, James Farwell, Sameh Gobriel, Charlie Tai, and Anshu-
368 mali Shrivastava. Slide: In defense of smart algorithms over hardware acceleration
369 for large-scale deep learning systems. In *In Proceedings of the 3rd Conference on*
370 *Machine Learning and Systems (MLSys)*, 2020.

371 [CT17] Timothy M Chan and Konstantinos Tsakalidis. Dynamic orthogonal range search-
372 ing on the ram, revisited. *Leibniz International Proceedings in Informatics, LIPIcs*,
373 77:281–2813, 2017.

374 [CWB⁺11] Ronan Collobert, Jason Weston, Léon Bottou, Michael Karlen, Koray Kavukcuoglu,
375 and Pavel Kuksa. Natural language processing (almost) from scratch. *Journal of*
376 *machine learning research*, 12(ARTICLE):2493–2537, 2011.

377 [DCLT18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-
378 training of deep bidirectional transformers for language understanding. *arXiv preprint*
379 *arXiv:1810.04805*, 2018.

380 [DIIM04] Mayur Datar, Nicole Immorlica, Piotr Indyk, and Vahab S Mirrokni. Locality-
381 sensitive hashing scheme based on p-stable distributions. In *Proceedings of the twen-*
382 *tieth annual symposium on Computational geometry (SoCG)*, pages 253–262, 2004.

383 [DIRW20] Yihe Dong, Piotr Indyk, Ilya Razenshteyn, and Tal Wagner. Learning space partitions
384 for nearest neighbor search. In *ICLR*. arXiv preprint arXiv:1901.08544, 2020.

385 [DLL⁺19] Simon S Du, Jason D Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient
386 descent finds global minima of deep neural networks. In *ICML*. [https://arxiv.](https://arxiv.org/pdf/1811.03804)
387 [org/pdf/1811.03804](https://arxiv.org/pdf/1811.03804), 2019.

388 [DMZS21] Shabnam Daghighi, Nicholas Meisburger, Mengnan Zhao, and Anshumali Shrivas-
389 tava. Accelerating slide deep learning on modern cpus: Vectorization, quantizations,
390 memory optimizations, and more. *Proceedings of Machine Learning and Systems*, 3,
391 2021.

392 [DZPS19] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably
393 optimizes over-parameterized neural networks. In *ICLR*, 2019.

394 [HLSY21] Baihe Huang, Xiaoxiao Li, Zhao Song, and Xin Yang. Fl-ntk: A neural tangent
395 kernel-based framework for federated learning convergence analysis. In *ICML*, 2021.

396 [Hoe63] Wassily Hoeffding. Probability inequalities for sums of bounded random variables.
397 *Journal of the American Statistical Association*, 58(301):13–30, 1963.

398 [HZRS16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning
399 for image recognition. In *Proceedings of the IEEE conference on computer vision and*
400 *pattern recognition (CVPR)*, pages 770–778, 2016.

401 [IM98] Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: towards removing
402 the curse of dimensionality. In *Proceedings of the thirtieth annual ACM symposium*
403 *on Theory of computing (STOC)*, pages 604–613, 1998.

404 [JKL⁺20] Haotian Jiang, Tarun Kathuria, Yin Tat Lee, Swati Padmanabhan, and Zhao Song. A
405 faster interior point method for semidefinite programming. In *FOCS*, 2020.

406 [JLSW20] Haotian Jiang, Yin Tat Lee, Zhao Song, and Sam Chiu-wai Wong. An improved
407 cutting plane method for convex optimization, convex-concave games and its appli-
408 cations. In *STOC*, 2020.

- 409 [KKL20] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient trans-
410 former. *arXiv preprint arXiv:2001.04451*, 2020.
- 411 [KSH12] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification
412 with deep convolutional neural networks. *Advances in neural information processing
413 systems*, 25:1097–1105, 2012.
- 414 [LBBH98] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based
415 learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–
416 2324, 1998.
- 417 [LL18] Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via
418 stochastic gradient descent on structured data. In *NeurIPS*, 2018.
- 419 [LS01] W.V. Li and Q.-M. Shao. Gaussian processes: Inequalities, small ball probabilities and
420 applications. In *Stochastic Processes: Theory and Methods*, volume 19 of *Handbook
421 of Statistics*, pages 533–597. Elsevier, 2001.
- 422 [LSS⁺20] Jason D Lee, Ruoqi Shen, Zhao Song, Mengdi Wang, and Zheng Yu. Generalized
423 leverage score sampling for neural networks. In *NeurIPS*, 2020.
- 424 [LSZ19] Yin Tat Lee, Zhao Song, and Qiuyi Zhang. Solving empirical risk minimization in the
425 current matrix multiplication time. In *Conference on Learning Theory (COLT)*, pages
426 2140–2157. PMLR, 2019.
- 427 [LXJ⁺20] Zichang Liu, Zhaozhuo Xu, Alan Ji, Jonathan Li, Beidi Chen, and Anshumali
428 Shrivastava. Climbing the wol: Training for cheaper inference. *arXiv preprint
429 arXiv:2007.01230*, 2020.
- 430 [Mat92a] Jiří Matoušek. Efficient partition trees. *Discrete & Computational Geometry*,
431 8(3):315–334, 1992.
- 432 [Mat92b] Jiri Matousek. Reporting points in halfspaces. *Computational Geometry*, 2(3):169–
433 186, 1992.
- 434 [OS20] Samet Oymak and Mahdi Soltanolkotabi. Towards moderate overparameterization:
435 global convergence guarantees for training shallow neural networks. In *arXiv preprint.
436 <https://arxiv.org/pdf/1902.04674.pdf>*, 2020.
- 437 [Raz17] Ilya Razenshteyn. *High-dimensional similarity search and sketching: algorithms and
438 hardness*. PhD thesis, Massachusetts Institute of Technology, 2017.
- 439 [Sch11] Jssai Schur. Bemerkungen zur theorie der beschränkten bilinearformen mit unendlich
440 vielen veränderlichen. 1911.
- 441 [SHM⁺16] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George
442 Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam,
443 Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree
444 search. *nature*, 529(7587):484–489, 2016.
- 445 [SL14] Anshumali Shrivastava and Ping Li. Asymmetric lsh (alsh) for sublinear time maxi-
446 mum inner product search (mips). *Advances in Neural Information Processing Sys-
447 tems (NIPS)*, pages 2321–2329, 2014.
- 448 [SL15a] Anshumali Shrivastava and Ping Li. Asymmetric minwise hashing for indexing bi-
449 nary inner products and set containment. In *Proceedings of the 24th international
450 conference on world wide web (WWW)*, pages 981–991, 2015.
- 451 [SL15b] Anshumali Shrivastava and Ping Li. Improved asymmetric locality sensitive hashing
452 (alsh) for maximum inner product search (mips). In *Proceedings of the Thirty-First
453 Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 812–821, 2015.

- 454 [SLJ⁺15] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir
455 Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going
456 deeper with convolutions. In *Proceedings of the IEEE conference on computer vi-*
457 *sion and pattern recognition*, pages 1–9, 2015.
- 458 [SSS⁺17] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang,
459 Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mas-
460 tering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- 461 [SY19] Zhao Song and Xin Yang. Quadratic suffices for over-parametrization via matrix
462 chernoff bound. *arXiv preprint arXiv:1906.03593*, 2019.
- 463 [TOG17] Csaba D Toth, Joseph O’Rourke, and Jacob E Goodman. *Handbook of discrete and*
464 *computational geometry*. CRC press, 2017.
- 465 [Tro15] Joel A Tropp. An introduction to matrix concentration inequalities. *Foundations and*
466 *Trends in Machine Learning*, 8(1-2):1–230, 2015.
- 467 [Vai89] Pravin M Vaidya. Speeding-up linear programming using fast matrix multiplication.
468 In *FOCS*, 1989.
- 469 [ZPD⁺20] Yi Zhang, Orestis Plevrakis, Simon S Du, Xingguo Li, Zhao Song, and Sanjeev Arora.
470 Over-parameterized adversarial training: An analysis overcoming the curse of dimen-
471 sionality. In *NeurIPS*. arXiv preprint arXiv:2002.06668, 2020.

472 Checklist

- 473 1. For all authors...
- 474 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
475 contributions and scope?
476 A: [Yes] .
- 477 (b) Did you describe the limitations of your work?
478 A: [Yes] See Section 7.
- 479 (c) Did you discuss any potential negative societal impacts of your work?
480 A: [N/A] . Our theoretical work does not have explicitly negative societal impacts.
- 481 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
482 them?
483 A: [Yes] , we conformed.
- 484 2. If you are including theoretical results...
- 485 (a) Did you state the full set of assumptions of all theoretical results?
486 A: [Yes] , we explicitly stated the assumptions.
- 487 (b) Did you include complete proofs of all theoretical results?
488 A: [Yes] , we provided the complete proofs in supplementary materials.
- 489 3. If you ran experiments...
- 490 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
491 mental results (either in the supplemental material or as a URL)?
492 A: [N/A]
- 493 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
494 were chosen)?
495 A: [N/A]
- 496 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
497 ments multiple times)?
498 A: [N/A]
- 499 (d) Did you include the total amount of compute and the type of resources used (e.g., type
500 of GPUs, internal cluster, or cloud provider)?
501 A: [N/A]
- 502 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

503 (a) If your work uses existing assets, did you cite the creators?
504 A: [N/A]

505 (b) Did you mention the license of the assets?
506 A: [N/A]

507 (c) Did you include any new assets either in the supplemental material or as a URL?
508 A: [N/A]

509 (d) Did you discuss whether and how consent was obtained from people whose data
510 you're using/curating?
511 A: [N/A]

512 (e) Did you discuss whether the data you are using/curating contains personally identifi-
513 able information or offensive content?
514 A: [N/A]

515 5. If you used crowdsourcing or conducted research with human subjects...

516 (a) Did you include the full text of instructions given to participants and screenshots, if
517 applicable?
518 A: [N/A]

519 (b) Did you describe any potential participant risks, with links to Institutional Review
520 Board (IRB) approvals, if applicable?
521 A: [N/A]

522 (c) Did you include the estimated hourly wage paid to participants and the total amount
523 spent on participant compensation?
524 A: [N/A]