

# End-to-end Learning of a Convolutional Neural Network via Deep Tensor Decomposition

Samet Oymak\* and Mahdi Soltanolkotabi†

May 18, 2018

## Abstract

In this paper we study the problem of learning the weights of a deep convolutional neural network. We consider a network where convolutions are carried out over non-overlapping patches with a single kernel in each layer. We develop an algorithm for simultaneously learning all the kernels from the training data. Our approach dubbed Deep Tensor Decomposition (DeepTD<sup>1</sup>) is based on a rank-1 tensor decomposition. We theoretically investigate DeepTD under a realizable model for the training data where the inputs are chosen i.i.d. from a Gaussian distribution and the labels are generated according to planted convolutional kernels. We show that DeepTD is data-efficient and provably works as soon as the sample size exceeds the total number of convolutional weights in the network. We carry out a variety of numerical experiments to investigate the effectiveness of DeepTD and verify our theoretical findings.

## 1 Introduction

Deep neural network (DNN) architectures have led to state of the art performance in many domains including image recognition, natural language processing, recommendation systems, and video analysis [9, 16, 18, 22, 39]. Convolutional neural networks (CNNs) are a class of deep, feed-forward neural networks with a specialized DNN architecture. CNNs are responsible for some of the most significant performance gains of DNN architectures. In particular, CNN architectures have led to striking performance improvements for image/object recognition tasks. Convolutional neural networks, loosely inspired by the visual cortex of animals, construct increasingly higher level features (such as mouth and nose) from lower level features such as pixels. An added advantage of CNNs which makes them extremely attractive for large-scale applications is their remarkable efficiency which can be attributed to: (1) intelligent utilization of parameters via weight-sharing, (2) their convolutional nature which exploits the local spatial structure of images/videos effectively, and (3) highly efficient matrix/vector multiplication involved in CNNs compared to fully-connected neural network architectures.

Despite the wide empirical success of CNNs the reasons for the effectiveness of neural networks and CNNs in particular is still a mystery. Recently there has been a surge of interest in developing more rigorous foundations for neural networks [1, 19, 23, 25, 27, 32, 34, 40, 41]. Most of this existing literature however focus on learning shallow neural networks typically consisting of zero or one hidden layer. In practical applications, depth seems to play a crucial role in constructing progressively higher-level features from pixels. Indeed, state of the art Resnet models typically have hundreds of layers. Furthermore, recent results suggest that increasing depth may substantially boost the expressive power of neural networks [7, 29].

In this paper, we propose an algorithm for approximately learning an arbitrarily deep CNN model with rigorous guarantees. Our goal is to provide theoretical insights towards better understanding when training deep CNN architectures is computationally tractable and how much data is required for successful training. We focus on a realizable model where the inputs are chosen i.i.d. from a Gaussian distribution and the labels

---

\*Department of Electrical and Computer Engineering, University of California, Riverside, CA

†Ming Hsieh Department of Electrical Engineering, University of Southern California, Los Angeles, CA

<sup>1</sup>Pun intended.

are generated according to planted convolutional kernels. We use both labels and features in the training data to construct a tensor. Our first insight is that, in the limit of infinite data this tensor converges to a *population* tensor which is approximately rank one and whose factors reveal the direction of the kernels. Our second insight is that even with finite data this *empirical* tensor is still approximately rank one. We show that the gap between the population and empirical tensors provably decreases with the increase in the size of the training data set and becomes negligible as soon as the size of the training data becomes proportional to the total numbers of the parameters in the planted CNN model. Combining these insights we provide a tensor decomposition algorithm to learn the kernels from training data. We show that our algorithm approximately learns the kernels (up to sign/scale ambiguities) as soon as the size of the training data is proportional to the total number of parameters of the planted CNN model. Our results can be viewed as a first step towards provable end-to-end learning of practical deep CNN models. Extending the connections between neural networks and tensors [7, 19, 40], we show how tensor decomposition can be utilized to approximately learn deep networks despite the presence of nonlinearities and growing depth. While our focus in this work is limited to tensors, we believe that our proposed algorithm may provide valuable insights for initializing local search methods (such as stochastic gradient descent) to enhance the quality and/or speed of CNN training.

## 2 Problem formulation and models

In this section we discuss the CNN model which is the focus of this paper. Our exposition is directly adapted from the companion paper [28]. A fully connected artificial neural network is composed of computational units called neurons. The neurons are decomposed into layers consisting of one input layer, one output layer and a few hidden layers with the output of each layer is fed in (as input) to the next layer. In a CNN model the output of each layer is related to the input of the next layer by a convolution operation. In this paper we focus on CNN model where the stride length is equal to the length of the kernel. This is sometimes referred to as a non-overlapping convolution operation. We now formally define this operation for future reference.

**Definition 2.1 (Non-overlapping convolution)** For two vectors  $\mathbf{k} \in \mathbb{R}^d$  and  $\mathbf{h} \in \mathbb{R}^{p=d\bar{p}}$  their non-overlapping convolution, denoted by  $\mathbf{k} \boxtimes \mathbf{h}$  yields a vector  $\mathbf{u} \in \mathbb{R}^{\bar{p}=\frac{p}{d}}$  whose entries are given by

$$\mathbf{u}_i = \langle \mathbf{k}^{(\ell)}, \mathbf{h}[i] \rangle \quad \text{where} \quad \mathbf{h}[i] := \begin{bmatrix} \mathbf{h}_{(i-1)d+1} \\ \mathbf{h}_{(i-1)d+2} \\ \vdots \\ \mathbf{h}_{id} \end{bmatrix}.$$

As mentioned earlier the non-overlapping convolution resembles the standard convolution between two vectors  $\mathbf{k} \in \mathbb{R}^d$  and  $\mathbf{h} \in \mathbb{R}^p$ , denoted by  $\mathbf{k} \otimes \mathbf{h}$ , which yields a vector in  $\mathbb{R}^{d+p}$ . The main difference is that rather than sliding the kernel one entry at a time across overlapping patches of  $\mathbf{h}$  in the convolution summation we use a stride length equal to the size of the kernel (i.e. slide the kernel by the length of the kernel) so that the summation is carried out over non-overlapping patches of  $\mathbf{h}$ . In this paper it is often convenient to view convolutions as matrix/vector multiplications. This leads us to the definition of the kernel matrix below.

**Definition 2.2 (Kernel matrix)** Consider a kernel  $\mathbf{k} \in \mathbb{R}^d$  and any vector  $\mathbf{h} \in \mathbb{R}^{p=d\bar{p}}$ . Corresponding to the non-overlapping convolution  $\mathbf{k} \boxtimes \mathbf{h}$ , we associate a kernel matrix  $\mathbf{K}_{\boxtimes} \in \mathbb{R}^{\bar{p} \times p}$  defined as

$$\mathbf{K}_{\boxtimes} := \mathbf{I}_{\bar{p}} \otimes \mathbf{k}^T = \begin{bmatrix} \mathbf{k}^T & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{k}^T & \mathbf{0} & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{k}^T \end{bmatrix}.$$

Here,  $\mathbf{A} \otimes \mathbf{B}$  denotes the Kronecker product between the two matrices  $\mathbf{A}$  and  $\mathbf{B}$  and  $\mathbf{I}_{\bar{p}}$  denotes the  $\bar{p} \times \bar{p}$  identity matrix. We note that based on this definition  $\mathbf{k} \boxtimes \mathbf{h} = \mathbf{K}_{\boxtimes} \mathbf{h}$ . Throughout the paper we shall use  $\mathbf{K}$  interchangeably with  $\mathbf{K}_{\boxtimes}$  to denote this kernel matrix with the dependence on the underlying kernel and its non-overlapping form implied.

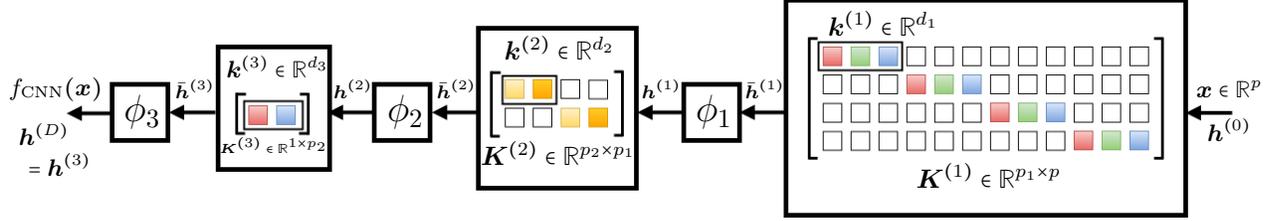


Figure 1: Depiction of the input-output relationship of a non-overlapping Convolutional Neural Network (CNN) model along with the various notations and symbols.

With the definition of the non-overlapping convolution and the corresponding kernel matrix in hand we are now ready to define the CNN model which is the focus of this paper. For the convenience of the reader the CNN input-output relationship along with the corresponding notation is depicted in Figure 1.

- **Depth and numbering of the layers.** We consider a network of depth  $D$  where we number the input as layer 0 and the output as layer  $D$  and the hidden layers 1 to  $D - 1$ .
- **Layer dimensions and representations.** We assume the input of the CNN, denoted by  $\mathbf{x} \in \mathbb{R}^p$ , consists of  $p$  features and the output is a one dimensional label. We also assume the hidden layers (numbered by  $\ell = 1, 2, \dots, D - 1$ ) consists of  $p_\ell$  units with  $\bar{\mathbf{h}}^{(\ell)} \in \mathbb{R}^{p_\ell}$  and  $\mathbf{h}^{(\ell)} \in \mathbb{R}^{p_\ell}$  denoting the input and output values of the units in the  $\ell$ th hidden layer. For consistency of our notation we shall also define  $\mathbf{h}^{(0)} := \mathbf{x} \in \mathbb{R}^p$  and note that the output of the CNN is  $\mathbf{h}^{(D)} \in \mathbb{R}$ . Furthermore,  $p_0 = p$  and  $p_D = 1$ .
- **Kernel dimensions and representation.** For  $\ell = 1, \dots, D$  we assume the kernel relating the output of layer  $(\ell - 1)$  to the input of layer  $\ell$  is of dimension  $d_\ell$  and is denoted by  $\mathbf{k}^{(\ell)} \in \mathbb{R}^{d_\ell}$ .
- **Inter-layer relationship.** We assume the inputs of layer  $\ell$  (denoted by  $\bar{\mathbf{h}}^{(\ell)} \in \mathbb{R}^{p_\ell}$ ) are related to the outputs of layer  $(\ell - 1)$  (denoted by  $\mathbf{h}^{(\ell-1)} \in \mathbb{R}^{p_{\ell-1}}$ ) via a non-overlapping convolution

$$\bar{\mathbf{h}}^{(\ell)} = \mathbf{k}^{(\ell)} \boxtimes \mathbf{h}^{(\ell-1)} = \mathbf{K}^{(\ell)} \mathbf{h}^{(\ell-1)} \quad \text{for } \ell = 1, \dots, D.$$

In the latter equality we have used the representation of non-overlapping convolution as a matrix/vector product involving the kernel matrix  $\mathbf{K}^{(\ell)} \in \mathbb{R}^{p_\ell \times p_{\ell-1}}$  associated with the kernel  $\mathbf{k}^{(\ell)} \in \mathbb{R}^{d_\ell}$  per Definition 2.2. We note that by the nature of the non-overlapping convolution we have that  $p_\ell = p_{\ell-1}/d_\ell$ .

- **Activation functions and intra-layer relationship.** We assume the input of each hidden unit is related to its output by applying an activation function  $\phi_\ell : \mathbb{R} \rightarrow \mathbb{R}$ . More precisely,  $\mathbf{h}^{(\ell)} := \phi_\ell(\bar{\mathbf{h}}^{(\ell)})$  where for a vector  $\mathbf{u} \in \mathbb{R}^p$ ,  $\phi_\ell(\mathbf{u}) \in \mathbb{R}^p$  is a vector obtained by applying the activation function  $\phi_\ell$  to each of the entries of  $\mathbf{u}$ . We allow for using distinct activation functions  $\{\phi_\ell\}_{\ell=1}^D$  at every layer. Throughout, we also assume all activations are 1-Lipschitz functions (i.e.  $|\phi_\ell(a) - \phi_\ell(b)| \leq |a - b|$ ).
- **Final output.** To summarize the input-output relation of our CNN model with an input  $\mathbf{x} \in \mathbb{R}^p$  and input kernel  $\mathbf{k}^{(1)} \in \mathbb{R}^{d_1}$  is given by

$$\mathbf{x} \mapsto f_{CNN}(\mathbf{x}) := \mathbf{h}^{(D)}, \quad (2.1)$$

with

$$\mathbf{h}^{(\ell)} = \phi_\ell(\bar{\mathbf{h}}^{(\ell)}) \quad \text{and} \quad \bar{\mathbf{h}}^{(\ell)} = \mathbf{K}^{(\ell)} \mathbf{h}^{(\ell-1)} = \mathbf{K}^{(\ell)} \phi_{\ell-1} \left( \mathbf{K}^{(\ell-1)} \phi_{\ell-2} \left( \dots \left( \mathbf{K}^{(2)} \phi_1 \left( \mathbf{K}^{(1)} \mathbf{x} \right) \right) \right) \right). \quad (2.2)$$

### 3 Algorithm: Deep Tensor Decomposition (DeepTD)

This paper introduces an approach to approximating the convolutional kernels from training data based on tensor decompositions dubbed DeepTD, which consists of a carefully designed tensor decomposition. To connect these two problems, we begin by stating how we intend to construct the tensor from the training data. To this aim given any input data  $\mathbf{x} \in \mathbb{R}^p$ , we form a  $D$ -way tensor  $\mathbf{X} \in \mathbb{R}^{\otimes_{\ell=1}^D d_\ell}$  as follows. First, we convert  $\mathbf{x}$  into a matrix by placing every  $d_1$  consecutive entries of  $\mathbf{x}$  as a row of a matrix of size  $p_1 \times d_1$ . From this matrix we then create a 3-way tensor of size  $p_2 \times d_2 \times d_1$  by grouping  $d_2$  consecutive entries of each of the  $d_1$  columns and so on. We repeat this procedure  $D$  times to arrive at the  $D$ -way tensor  $\mathbf{X} \in \mathbb{R}^{\otimes_{\ell=1}^D d_\ell}$ . We define  $\mathcal{T} : \mathbb{R}^p \mapsto \mathbb{R}^{\otimes_{\ell=1}^D d_\ell}$  as the corresponding tensor operation that maps  $\mathbf{x} \in \mathbb{R}^p$  to  $\mathbf{X} \in \mathbb{R}^{\otimes_{\ell=1}^D d_\ell}$ .

Given a set of training data consisting of  $n$  input/output pairs  $(\mathbf{x}_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$  we construct a tensor  $\mathbf{T}_n$  by tensorizing the input vectors as discussed above and calculating a weighted combination of these tensorized inputs where the weights depends on the output labels for the training data. More precisely,

$$\mathbf{T}_n := \frac{1}{n} \sum_{i=1}^n (y_i - y_{avg}) \mathbf{X}_i \quad \text{where} \quad y_{avg} = \frac{1}{n} \sum_{i=1}^n y_i \quad \text{and} \quad \mathbf{X}_i = \mathcal{T}(\mathbf{x}_i). \quad (3.1)$$

We then perform a rank-1 tensor decomposition on this  $\mathbf{T}_n$  tensor to approximate the convolutional kernels. Specifically we solve

$$\hat{\mathbf{k}}^{(1)}, \dots, \hat{\mathbf{k}}^{(D)} = \arg \max_{\mathbf{v}_1 \in \mathbb{R}^{d_1}, \mathbf{v}_2 \in \mathbb{R}^{d_2}, \dots, \mathbf{v}_D \in \mathbb{R}^{d_D}} \left\langle \mathbf{T}_n, \bigotimes_{\ell=1}^D \mathbf{v}_\ell \right\rangle \quad \text{subject to} \quad \|\mathbf{v}_1\|_{\ell_2} = \|\mathbf{v}_2\|_{\ell_2} = \dots = \|\mathbf{v}_D\|_{\ell_2} = 1. \quad (3.2)$$

In the above  $\bigotimes_{\ell=1}^D \mathbf{v}_\ell$  denotes the tensor resulting from the outer product of the vectors  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_D$ . This tensor rank decomposition is also known as CANDECOMP/PARAFAC (CP) decomposition [5] and can be solved efficiently using Alternating Least Squares (ALS) and a variety of other algorithms [2, 3, 14].<sup>2</sup>

At this point it is completely unclear why the tensor  $\mathbf{T}_n$  or its rank-1 decomposition can yield anything useful. The main intuition is that as the data set grows ( $n \rightarrow \infty$ ) the *empirical* tensor  $\mathbf{T}_n$  converges close to a *population* tensor  $\mathbf{T}$  whose rank-1 decomposition reveals useful information about the kernels. Specifically, we will show that

$$\lim_{n \rightarrow \infty} \mathbf{T}_n = \mathbf{T} := \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)} [f_{\text{CNN}}(\mathbf{x}) \mathcal{T}(\mathbf{x})] \approx \alpha \bigotimes_{\ell=1}^D \mathbf{k}^{(\ell)}, \quad (3.3)$$

with  $\alpha$  a scalar whose value shall be discussed later on. Here,  $\mathbf{x}$  is a Gaussian random vector with i.i.d.  $\mathcal{N}(0, 1)$  entries and represents a typical input with  $f_{\text{CNN}}(\mathbf{x})$  the corresponding output and  $\mathcal{T}(\mathbf{x})$  the tensorized input. We will also utilize a concentration argument to show that when the training data set originates from an i.i.d. distribution, for a sufficiently large training data  $n$ ,  $\mathbf{T}_n$  yields a good approximation of the population tensor  $\mathbf{T}$ .

Another perhaps perplexing aspect of the construction of  $\mathbf{T}_n$  in (3.1) is the subtraction by  $y_{avg}$  in the weights. The reason this may be a source of confusion is that based on the intuition above

$$\mathbb{E}[\mathbf{T}] = \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n y_i \mathbf{X}_i \right] = \mathbb{E} \left[ \frac{1}{n-1} \sum_{i=1}^n (y_i - y_{avg}) \mathbf{X}_i \right] = \frac{n}{n-1} \mathbb{E}[\mathbf{T}_n] \approx \alpha \bigotimes_{\ell=1}^D \mathbf{k}^{(\ell)},$$

so that the subtraction by the average seems completely redundant. The main purpose of this subtraction is to ensure the weights  $y_i - y_{avg}$  are centered (have mean zero). This centering allows for a much better concentration of the empirical tensor around its population counter part and is crucial to the success of our approach. We would like to point out that such a centering procedure is reminiscent of batch-normalization heuristics deployed when training deep neural networks.

<sup>2</sup>We would like to note that we believe that the assumptions under which the above stated algorithms solve the problem (3.2) holds for the empirical tensor  $\mathbf{T}_n$  under our planted model. We do not however pursue this here and leave it to future work.

Finally, we note that based on (3.3), the rank-1 tensor decomposition step can recover the convolutional kernels  $\{\mathbf{k}^{(\ell)}\}_{\ell=1}^D$  up to sign and scaling ambiguities. Unfortunately, depending on the activation function, it may be impossible to overcome these ambiguities. For instance, if the activations are homogeneous (i.e.  $\phi_\ell(ax) = a\phi_\ell(x)$ ), then scaling up one layer and scaling down the other layer by the same amount does not change the overall function  $f_{\text{CNN}}(\cdot)$ . Similarly, if the activations are odd functions, negating two of the layers at the same time preserves the overall function. In Section 5 we discuss some heuristics and theoretical guarantees for overcoming these sign/scale ambiguities.

## 4 Main results

In this section we introduce our theoretical results for DeepTD. We will discuss these results in three sections. In Section 4.1 we show that the empirical tensor concentrates around its population counterpart. Then in Section 4.2 we show that the population tensor is well-approximated by a rank-1 tensor whose factors reveal the convolutional kernels. Finally, in Section 4.3 we combine these results to show DeepTD can approximately learn the convolutional kernels up to sign/scale ambiguities.

### 4.1 Concentration of the empirical tensor

Our first result shows that the empirical tensor concentrates around the population tensor. We measure the quality of this concentration via the tensor spectral norm defined below.

**Definition 4.1** Let  $\mathcal{RO}$  be the set of rank-one tensors characterized by

$$\mathcal{RO} = \left\{ \mathbf{V} \in \mathbb{R}^{\otimes_{\ell=1}^D d_\ell} \mid \mathbf{V} = \bigotimes_{\ell=1}^D \mathbf{v}_\ell, \|\mathbf{v}_i\|_{\ell_2} \leq 1 \text{ for all } 1 \leq i \leq D \right\}.$$

The spectral norm of a tensor  $\mathbf{X} \in \mathbb{R}^{\otimes_{\ell=1}^D d_\ell}$  is given by the supremum,

$$\|\mathbf{T}\| = \sup_{\mathbf{V} \in \mathcal{RO}} \langle \mathbf{V}, \mathbf{T} \rangle.$$

**Theorem 4.2** Consider a CNN model  $\mathbf{x} \mapsto f_{\text{CNN}}(\mathbf{x})$  of the form (2.1) consisting of  $D \geq 2$  layers with convolutional kernels  $\mathbf{k}^{(1)}, \mathbf{k}^{(2)}, \dots, \mathbf{k}^{(D)}$  of lengths  $d_1, d_2, \dots, d_D$ . Let  $\mathbf{x} \in \mathbb{R}^p$  be a Gaussian random vector distributed as  $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$  with the corresponding labels  $y = f_{\text{CNN}}(\mathbf{x})$  generated by the CNN model and  $\mathbf{X} := \mathcal{T}(\mathbf{x})$  the corresponding tensorized input. Suppose the data set consists of  $n$  training samples where the feature vectors  $\mathbf{x}_i \in \mathbb{R}^p$  are distributed i.i.d.  $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$  with the corresponding labels  $y_i = f_{\text{CNN}}(\mathbf{x}_i)$  generated by the same CNN model and  $\mathbf{X}_i := \mathcal{T}(\mathbf{x}_i)$  the corresponding tensorized input. Then the empirical tensor  $\mathbf{T}_n$  and population tensor  $\mathbf{T}$  defined based on this dataset obey

$$\|\mathbf{T}_n - \mathbf{T}\| := \left\| \frac{1}{n} \sum_{i=1}^n (y_i - y_{\text{avg}}) \mathbf{X}_i - \mathbb{E}[y \mathbf{X}] \right\| \leq c \prod_{\ell=1}^D \|\mathbf{k}^{(\ell)}\|_{\ell_2} \frac{\sqrt{(\sum_{\ell=1}^D d_\ell) \log D + t}}{\sqrt{n}}, \quad (4.1)$$

with probability at least  $1 - 5e^{-\min(t^2, t\sqrt{n}, n)}$ , where  $c > 0$  is an absolute constant.

The theorem above shows that the empirical tensor approximates the population tensor with high probability. This theorem also shows that the quality of this approximation is proportional to  $\prod_{\ell=1}^D \|\mathbf{k}^{(\ell)}\|_{\ell_2}$ . This is natural as  $\prod_{\ell=1}^D \|\mathbf{k}^{(\ell)}\|_{\ell_2}$  is an upper-bound on the Lipschitz constant of the network and shows how much the CNN output fluctuates with changes in the input. The more fluctuations, the less concentrated the empirical tensor is, translating into a worse approximation guarantee. Furthermore, the quality of this approximation grows with the square root of the parameters in the model ( $\sum_{\ell=1}^D d_\ell$ ) and is inversely proportional to the square root of the number of samples ( $n$ ) which are typical scalings in statistical learning. We would also like to note that as we will see in the forthcoming sections, in many cases  $\|\mathbf{T}\|$  is roughly on the order of

$\prod_{\ell=1}^D \|\mathbf{k}^{(\ell)}\|_{\ell_2}$  so that (4.1) guarantees that  $\|\mathbf{T}_n - \mathbf{T}\|/\|\mathbf{T}\| \leq c' \sqrt{\frac{(\sum_{\ell=1}^D d_\ell) \log D}{n}}$ . Therefore, the relative error in the approximation is less than  $\epsilon$  as soon as the number of observations exceeds the number of parameters in the model by a logarithmic factor in depth i.e.  $n \gtrsim (\sum_{\ell=1}^D d_\ell) \frac{\log D}{\epsilon^2}$ .

## 4.2 Rank one approximation of the population tensor

Our second result shows that the population tensor can be approximated by a rank one tensor. To explain the structure of this rank one tensor and quantify the quality of this approximation we require a few definitions. The first quantity roughly captures the average amount by which the nonlinear activations amplify or attenuate the size of an input feature at the output. This quantity is formally defined below.

**Definition 4.3 (CNN gain)** *Let  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_p)$  and define the hidden unit/output values of the CNN based on this random input per equations (2.2) and (2.1). We define the CNN gain as*

$$\alpha_{CNN} = \prod_{\ell=1}^D \mathbb{E}[\phi'_\ell(\bar{\mathbf{h}}_1^{(\ell)})].$$

*In words, this is the product of expectations of the activations evaluated at the first entry of each layer.*

This quantity is the product of the average slopes of the activations evaluated along a path connecting the first input feature to the first hidden units across the layers all the way to the output. We note that this quantity is the same when calculated along any path connecting an input feature to the output passing through the hidden units. Therefore, this quantity can be thought of as the average gain (amplification or attenuation) of a given input feature due to the nonlinear activations in the network. To gain some intuition consider a ReLU network which is mostly inactive. Then the network is dead and  $\alpha_{CNN} \approx 0$ . On the other extreme if all ReLU units are active the network operates in the linear regime and  $\alpha_{CNN} = 1$ . We would like to point out that  $\alpha_{CNN}$  can in many cases be bounded from below by a constant. For instance for ReLU activations as long as the kernels obey

$$(\mathbf{1}^T \mathbf{k}^{(\ell)}) \geq 4 \|\mathbf{k}^{(\ell)}\|_{\ell_2}, \quad (4.2)$$

then  $\alpha_{CNN} \geq \gamma$  with  $\gamma$  a fixed numerical constant (formally proved in our companion paper [28]). We note that an assumption similar to (4.2) is needed for the network to be active. This is because if the kernel sums are negative one can show that with high probability, all the ReLUs after the first layer will be inactive and the network will be dead. In [28] we also show that similar results hold for other popular activations including softplus. With this definition in hand, we are now ready to describe the form of the rank one tensor that approximates the population tensor.

**Definition 4.4 (Rank one CNN tensor)** *We define the rank one CNN tensor  $\mathbf{L}_{CNN} \in \mathbb{R}^{\otimes_{\ell=1}^D d_\ell}$  as*

$$\mathbf{L}_{CNN} = \alpha_{CNN} \bigotimes_{\ell=1}^D \mathbf{k}^{(\ell)}.$$

*In words, this is the product of the kernels  $\{\mathbf{k}^{(\ell)}\}_{\ell=1}^D$  scaled by the CNN gain  $\alpha_{CNN}$ .*

To quantify how well the rank one CNN tensor approximates the population tensor we need two definitions. The first definition concerns the activation functions.

**Definition 4.5 (Activation smoothness)** *We assume the activations are differentiable everywhere and  $S$ -smooth (i.e.  $|\phi'_\ell(x) - \phi'_\ell(y)| \leq S|x - y|$  for all  $x, y \in \mathbb{R}$ ) for some  $S \geq 0$ .*

The reason smoothness of the activations play a role in the quality of the rank one approximation is that smoother activations translate into smoother variations in the entries of the population tensor. Therefore, the population tensor can be better approximated by a low-rank tensor. The second definition captures how diffused the kernels are.

**Definition 4.6 (Kernel diffuseness parameter)** Given kernels  $\{\mathbf{k}^{(\ell)}\}_{\ell=1}^D$  with dimensions  $\{d_\ell\}_{\ell=1}^D$ , the kernel diffuseness parameter  $\mu$  is defined as

$$\mu = \sup_{1 \leq \ell \leq D} \frac{\sqrt{d_\ell} \|\mathbf{k}^{(\ell)}\|_{\ell_\infty}}{\|\mathbf{k}^{(\ell)}\|_{\ell_2}}.$$

The less diffused (or more spiky) the kernels are, the more the population tensor fluctuates and thus the quality of the approximation to a rank one tensor decreases. With these definitions in place, we are now ready to state our theorem on approximating a population tensor with a rank one tensor.

**Theorem 4.7** Consider the setup of Theorem 4.2. Furthermore, assume the activations are  $S$ -smooth per Definition 4.5 and the convolutional kernels are  $\mu$ -diffused per Definition 4.6. Then, the population tensor  $\mathbf{T} := \mathbb{E}[y\mathbf{X}]$  can be approximated by the rank-1 tensor  $\mathbf{L}_{CNN} := \alpha_{CNN} \otimes_{\ell=1}^D \mathbf{k}^{(\ell)}$  as follows

$$\|\mathbf{T} - \mathbf{L}_{CNN}\| \leq \|\mathbf{T} - \mathbf{L}_{CNN}\|_F \leq \sqrt{8\pi} \mu S \cdot \prod_{i=1}^D \|\mathbf{k}^{(i)}\|_{\ell_2} \cdot \sup_{\ell} \prod_{i=1}^{\ell} \|\mathbf{k}^{(i)}\|_{\ell_2} \frac{D}{\sqrt{\min_{\ell} d_\ell}}.$$

The theorem above states that the quality of the rank one approximation deteriorates with increase in the smoothness of the activations and the diffuseness of the convolutional kernels. As mentioned earlier increase in these parameters leads to more fluctuations in the population tensor making it less likely that it can be well approximated by a rank one tensor. We also note that  $\|\mathbf{L}_{CNN}\| = \alpha_{CNN} \prod_{\ell=1}^D \|\mathbf{k}^{(\ell)}\|_{\ell_2}$  and therefore the relative error in this approximation is bounded by

$$\frac{\|\mathbf{T} - \mathbf{L}_{CNN}\|}{\|\mathbf{L}_{CNN}\|} \leq \sqrt{8\pi} \frac{\mu S}{\alpha_{CNN}} \sup_{\ell} \prod_{i=1}^{\ell} \|\mathbf{k}^{(i)}\|_{\ell_2} \frac{D}{\sqrt{\min_{\ell} d_\ell}}.$$

We would like to note that for many activations the smoothness is bounded by a constant. For instance, for the softplus activation ( $\phi(x) = \log(1 + e^x)$ ) and one can show that  $S \leq 1$ . As stated earlier, under appropriate assumptions on the kernels and activations, the CNN gain  $\alpha_{CNN}$  is also bounded from below by a constant. Assume the convolutional kernels have unit norm and are sufficiently diffused so that the diffuseness parameter is bounded by a constant. We can then conclude that

$$\frac{\|\mathbf{T} - \mathbf{L}_{CNN}\|}{\|\mathbf{L}_{CNN}\|} \leq \frac{c}{\alpha_{CNN}} \frac{D}{\sqrt{\min_{\ell} d_\ell}}.$$

This implies that as soon as the length of the convolutional patches scale with the square of depth of the network by a constant factor the rank one approximation is sufficiently good. Our back-of-the-envelope calculations suggest that the correct scaling is linear in  $D$  versus the quadratic result we have established here. Improving our result to achieve the correct scaling is an interesting future research direction. Finally, we would like to note that while we have assumed differentiable and smooth activations we expect our results to apply to popular non-differentiable activations such as ReLU activations. Indeed, as it will become clear in our proofs (e.g. see Lemma 9.14) an average notion of smoothness appears to be sufficient for our results to apply.

### 4.3 Learning the convolutional kernels

We demonstrated in the previous two sections that the empirical tensor concentrates around its population counter part and that the population tensor is well-approximated by a rank one tensor. We combine these two results along with a perturbation argument to provide guarantees for the DeepTD algorithm.

**Theorem 4.8 (Main theorem)** Consider a CNN model  $\mathbf{x} \mapsto f_{CNN}(\mathbf{x})$  of the form (2.1) consisting of  $D \geq 2$  layers with convolutional kernels  $\mathbf{k}^{(1)}, \mathbf{k}^{(2)}, \dots, \mathbf{k}^{(D)}$  of lengths  $d_1, d_2, \dots, d_D$ . Let  $\mathbf{x} \in \mathbb{R}^p$  be a Gaussian

random vector distributed as  $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$  with the corresponding labels  $y = f_{\text{CNN}}(\mathbf{x})$  generated by the CNN model and  $\mathbf{X} := \mathcal{T}(\mathbf{x})$  the corresponding tensorized input. Suppose the data set consists of  $n$  training samples where the feature vectors  $\mathbf{x}_i \in \mathbb{R}^p$  are distributed i.i.d.  $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$  with the corresponding labels  $y_i = f_{\text{CNN}}(\mathbf{x}_i)$  generated by the same CNN model and  $\mathbf{X}_i := \mathcal{T}(\mathbf{x}_i)$  the corresponding tensorized input. Also, assume the activations are 1-Lipschitz (i.e.  $|\phi'_\ell(z)| \leq 1$ ), are  $S$ -smooth per Definition 4.5 and the convolutional kernels are  $\mu$ -diffused per Definition 4.6. Then the empirical tensor  $\mathbf{T}_n := \frac{1}{n} \sum_{i=1}^n (y_i - y_{\text{avg}}) \mathbf{X}_i$  and the rank one approximation  $\mathbf{L}_{\text{CNN}} = \alpha_{\text{CNN}} \otimes_{\ell=1}^D \mathbf{k}^{(\ell)}$  to the population tensor  $\mathbf{T} = \mathbb{E}[\mathbf{y}\mathbf{X}]$  obey

$$\|\mathbf{T}_n - \mathbf{L}_{\text{CNN}}\| \leq \prod_{\ell=1}^D \|\mathbf{k}^{(\ell)}\|_{\ell_2} \left( c \frac{\sqrt{(\sum_{\ell=1}^D d_\ell) \log D + t}}{\sqrt{n}} + \sqrt{8\pi} \mu S \cdot \sup_{\ell} \prod_{i=1}^{\ell} \|\mathbf{k}^{(i)}\|_{\ell_2} \frac{D}{\sqrt{\min_{\ell} d_\ell}} \right),$$

with probability at least  $1 - 5e^{-\min(t^2, t\sqrt{n}, n)}$ , where  $c > 0$  is an absolute constant. Furthermore, the DeepTD estimates of the convolutional kernels given by (3.2) using the empirical tensor  $\mathbf{T}_n$  obeys

$$\frac{\prod_{\ell=1}^D |\langle \mathbf{k}^{(\ell)}, \hat{\mathbf{k}}^{(\ell)} \rangle|}{\prod_{\ell=1}^D \|\mathbf{k}^{(\ell)}\|_{\ell_2}} \geq 1 - \frac{2}{\alpha_{\text{CNN}}} \left( c \frac{\sqrt{(\sum_{\ell=1}^D d_\ell) \log D + t}}{\sqrt{n}} + \sqrt{8\pi} \mu S \cdot \sup_{\ell} \prod_{i=1}^{\ell} \|\mathbf{k}^{(i)}\|_{\ell_2} \frac{D}{\sqrt{\min_{\ell} d_\ell}} \right),$$

with probability at least  $1 - 5e^{-\min(t^2, t\sqrt{n}, n)}$ , where  $c > 0$  is an absolute constant.

The above theorem is our main result on learning a non-overlapping CNN with a single kernel at each layer. It demonstrates that estimates  $\hat{\mathbf{k}}^{(\ell)}$  obtained by DeepTD have significant inner product with the ground truth kernels  $\mathbf{k}^{(\ell)}$  with high probability, using only few samples. Indeed, similar to the discussion after Theorem 4.7 assuming the activations are sufficiently smooth and the convolutional kernels are unit norm and sufficiently diffused, the theorem above can be simplified as follows

$$\frac{\prod_{\ell=1}^D |\langle \mathbf{k}^{(\ell)}, \hat{\mathbf{k}}^{(\ell)} \rangle|}{\prod_{\ell=1}^D \|\mathbf{k}^{(\ell)}\|_{\ell_2}} \geq 1 - c \left( \frac{\sqrt{(\sum_{\ell=1}^D d_\ell) \log D}}{\sqrt{n}} + \frac{D}{\sqrt{\min_{\ell} d_\ell}} \right).$$

Thus the kernel estimates obtained via DeepTD are well aligned with the true kernels as soon as the number of samples scales with the total number of parameters in the model and the length of the convolutional kernels (i.e. the size of the batches) scales quadratically with the depth of the network.

## 5 Resolving sign and scaling ambiguities

We note that DeepTD operates by accurately approximating the rank one tensor  $\otimes_{\ell=1}^D \mathbf{k}^{(\ell)}$  from data. Therefore, DeepTD can only recover the convolutional kernels up to Sign/Scale Ambiguities (SSA). In general, it may not be possible to recover the ground truth kernels from the training data. For instance, when activations are ReLU, the norms of the kernels cannot be estimated from data as multiplying a kernel and dividing another by the same positive scalar leads to the same training data. However, we can try to learn a good approximation  $\hat{f}_{\text{CNN}}(\cdot)$  of the network  $f_{\text{CNN}}(\cdot)$  to minimize the risk  $\mathbb{E}[(f_{\text{CNN}}(\mathbf{x}) - \hat{f}_{\text{CNN}}(\mathbf{x}))^2]$ .

To this aim, we introduce Centered Empirical Risk Minimization (CERM) which is a slight modification of Empirical Risk Minimization (ERM). Let us first describe how finding a good  $\hat{f}_{\text{CNN}}(\cdot)$  can be formulated with CERM. Given  $n$  i.i.d. data points  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim (\mathbf{x}, y)$ , and a function class  $\mathcal{F}$ , CERM applies ERM after centering the residuals. Given  $f \in \mathcal{F}$ , define the average residual function  $r_{\text{avg}}(f) = \frac{1}{n} \sum_{i=1}^n y_i - f(\mathbf{x}_i)$ . We define the Centered Empirical Risk Minimizer as

$$\begin{aligned} \hat{f} &= \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i) - r_{\text{avg}}(f))^2, \\ &= \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i) - \mathbb{E}[(y_i - f(\mathbf{x}_i))])^2 - \frac{1}{n^2} \left( \sum_{i=1}^n y_i - f(\mathbf{x}_i) - \mathbb{E}[(y_i - f(\mathbf{x}_i))] \right)^2. \end{aligned} \quad (5.1)$$

The remarkable benefit of CERM over ERM is the fact that, the learning rate doesn't suffer from the label or function bias. This is in similar nature to the DeepTD algorithm that applies label centering. In the proofs (in particular Section 9.5, Theorem 9.17) we provide a generalization bound on the CERM solution (5.1) in terms of the Lipschitz covering number of the function space. While (5.1) can be used to learn all kernels, it does not provide an efficient algorithm. Instead, we will use CERM to resolve SSA after estimating the kernels via DeepTD. Interestingly, this approach only requires a few ( $\mathcal{O}(D)$ ) extra training samples. Inspired from CERM, in Section 5.1, we propose a greedy algorithm to address SSA. We will apply CERM to the following function class with bounded kernels,

$$\mathcal{F}_{\hat{\mathbf{k}}, B} := \left\{ f : \mathbb{R}^p \mapsto \mathbb{R} \mid f \text{ is a CNN function of the form (2.1) with kernels } \{\beta_\ell \hat{\mathbf{k}}^{(\ell)}\}_{\ell=1}^D \text{ with } |\beta_\ell| \leq B^{\frac{1}{D}} \right\}. \quad (5.2)$$

In words this is the function class of all CNN functions with kernels the same as those obtained by DeepTD up to sign/scale ambiguities  $\{\beta_\ell\}_{\ell=1}^D$  where the maximum scale ambiguity is  $B$ .

**Theorem 5.1** *Let  $f_{\text{CNN}}()$  be defined via (2.1) with convolutional kernels  $\{\mathbf{k}^{(\ell)}\}_{\ell=1}^D$  obeying  $\|\mathbf{k}^{(\ell)}\|_{\ell_2} \leq B^{1/D}$  for some  $B > 0$  and consider the function class  $\mathcal{F}_{\hat{\mathbf{k}}, B}$  above with the same choice of  $B$ . Assume we have  $n$  i.i.d. samples  $(\mathbf{x}_i, y_i) \sim (\mathbf{x}, y)$  where  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_p)$  and  $y = f_{\text{CNN}}(\mathbf{x})$ . Suppose for some  $\varepsilon \leq B$ ,*

$$n \geq cB^2 D \log\left(\frac{CDBp}{\varepsilon}\right) \max\left(\frac{1}{\varepsilon}, \frac{1}{\varepsilon^2}\right),$$

*holds for fixed numerical constants  $c, C > 0$ . Then the solution  $\hat{f}$  to the CERM problem (5.1) obeys*

$$\mathbb{E}\left[\left(\hat{f}(\mathbf{x}) - f_{\text{CNN}}(\mathbf{x})\right)^2\right] \leq \min_{f \in \mathcal{F}_{\hat{\mathbf{k}}, B}} \mathbb{E}[(f(\mathbf{x}) - f_{\text{CNN}}(\mathbf{x}))^2] + \varepsilon. \quad (5.3)$$

*on a new sample  $\mathbf{x} \in \mathcal{N}(0, \mathbf{I}_p)$  with probability at least  $1 - e^{-\gamma n} - 4n \exp(-p)$  with  $\gamma > 0$  an absolute constant.*

The above theorem states that CERM finds the sign/scale ambiguity that accurately estimates the labels on new data as long as the number of samples which are used in CERM exceeds the depth of the network by constant/log factors. In the next section we present a greedy heuristic for finding the CERM estimate.

## 5.1 Greedy algorithm for resolving sign and scale ambiguities

In order to resolve SSA, inspired from CERM, we propose Algorithm 1 which operates over the function class,

$$\mathcal{F}_{\hat{\mathbf{k}}} := \left\{ \gamma f : \mathbb{R}^p \mapsto \mathbb{R} \mid f \text{ is a CNN of the form (2.1) with kernels } \{\beta_\ell \hat{\mathbf{k}}^{(\ell)}\}_{\ell=1}^D \text{ with } \beta_\ell \in \{1, -1\}, \gamma \geq 0 \right\}. \quad (5.4)$$

It first determines the signs  $\beta_\ell$  by locally optimizing the kernels and then finds a global scaling  $\gamma > 0$ . In the first phase, the algorithm attempts to maximize the correlation between the centered labels  $y_{c,i} = y_i - n^{-1} \sum_{i=1}^n y_i$  and the  $\hat{f}_{\text{CNN}}()$  predictions given by  $\hat{y}_{c,i} = \hat{y}_i - n^{-1} \sum_{i=1}^n \hat{y}_i$ . It goes over all kernels one by one and it flips a kernel ( $\hat{\mathbf{k}}^{(\ell)} \rightarrow -\hat{\mathbf{k}}^{(\ell)}$ ) if flipping increases the correlation. This process goes on as long as there is an improvement. Afterwards, we use a simple linear regression to get the best scaling  $\gamma$  by minimizing the centered empirical loss  $\sum_{i=1}^n (y_{c,i} - \gamma \hat{y}_{c,i})^2$ . While our approach is applicable to arbitrary activations, it is tailored towards homogeneous activations ( $\phi(cx) = c\phi(x)$ ). The reason is that for homogeneous activations, function classes (5.2) and (5.4) coincide and a single global scaling  $\gamma$  is sufficient. Note that ReLU and the identity activation (i.e. no activation) are both homogeneous, in fact they are elements of a larger homogeneous activation family named Leaky ReLU. Leaky ReLU is parametrized by some scalar  $0 \leq \beta \leq 1$  and defined as follows

$$\text{LReLU}(x) = \begin{cases} x & \text{if } x \geq 0, \\ \beta x & \text{if } x < 0. \end{cases}$$

As we shall discuss in our numerical experiments in Section 6, Algorithm 1 works well for different depths ( $D = 4, 8, 12$ ) and kernel widths  $2 \leq d_\ell \leq 10$ . We would like to point out however that it appears that this local search heuristic can indeed get stuck in a suboptimal local minima even for a depth 4 convolutional network (see Figure 3 and the associated discussion). Developing efficient algorithms that can reliably construct a good  $\hat{f}_{\text{CNN}}()$  from  $\{\hat{\mathbf{k}}^{(\ell)}\}_{\ell=1}^D$  is an interesting direction for future research.

---

**Algorithm 1** Greedily algorithm for resolving sign/scale ambiguities for Leaky ReLU activations.

---

```

1: procedure MAXCORR
2: Inputs: Data  $(y_i, \mathbf{x}_i)_{i=1}^n$ , estimates  $\{\hat{\mathbf{k}}^{(\ell)}\}_{\ell=1}^D$ .
3:    $\rho_{\max} \leftarrow |\text{Corr}(\{\hat{\mathbf{k}}^{(\ell)}\}_{\ell=1}^D, 0)|$ , FLIP  $\leftarrow$  TRUE.
4:   while FLIP do
5:     FLIP  $\leftarrow$  FALSE.
6:     for  $1 \leq \ell \leq D$  do
7:        $\rho \leftarrow |\text{Corr}(\{\hat{\mathbf{k}}^{(1)}, \dots, -\hat{\mathbf{k}}^{(\ell)}, \dots, \hat{\mathbf{k}}^{(D)}\}, 0)|$ .
8:       if  $\rho > \rho_{\max}$  then
9:          $\rho_{\max} \leftarrow \rho$ 
10:         $\hat{\mathbf{k}}^{(\ell)} \leftarrow -\hat{\mathbf{k}}^{(\ell)}$ 
11:        FLIP  $\leftarrow$  TRUE
12:    $\gamma \leftarrow \text{Corr}(\{\hat{\mathbf{k}}^{(\ell)}\}_{\ell=1}^D, 1)$ .
13: return kernels  $\{\hat{\mathbf{k}}^{(\ell)}\}_{\ell=1}^D$ , scaling  $\gamma$ .
```

---

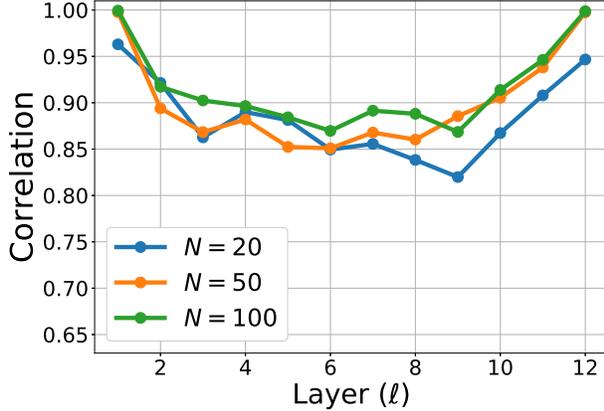
**Algorithm 2** Return the correlation between centered labels.

---

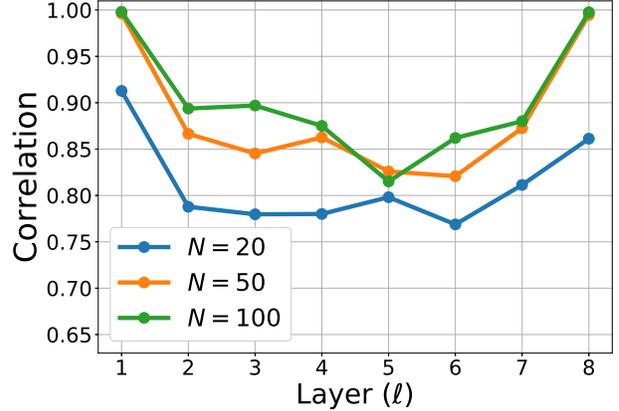
```

1: procedure CORR( $\{\hat{\mathbf{k}}^{(\ell)}\}_{\ell=1}^D$ , opt)
2:    $\hat{y}_i \leftarrow f_{\text{CNN}}(\{\hat{\mathbf{k}}^{(\ell)}\}_{\ell=1}^D; \mathbf{x}_i)$ .
3:    $y_{c,i} \leftarrow y_i - \frac{1}{n} \sum_{i=1}^n y_i$  and  $\hat{y}_{c,i} \leftarrow \hat{y}_i - \frac{1}{n} \sum_{i=1}^n \hat{y}_i$ .
4:    $\rho \leftarrow \sum_{i=1}^n y_{c,i} \hat{y}_{c,i}$ .
5: return  $\rho$  if opt = 0,  $\rho / (\sum_{i=1}^n \hat{y}_{c,i}^2)$  if opt = 1.
```

---



(a)  $d = 2, D = 12$ .



(b)  $d = 3, D = 8$ .

Figure 2: Correlations  $(\text{corr}(\hat{\mathbf{k}}^{(\ell)}, \mathbf{k}^{(\ell)}) = |\langle \hat{\mathbf{k}}^{(\ell)}, \mathbf{k}^{(\ell)} \rangle|)$  between the DeepTD estimate and the ground truth kernels for different layers and various over-sampling ratios ( $N = \frac{n}{\sum_{\ell=1}^D d_\ell}$ ).

## 6 Numerical experiments

Our goal in this section is to numerically corroborate the theoretical predictions of Section 4. To this aim we use a CNN model of the form (2.1) with  $D$  layers and ReLU activations and set the kernel lengths to be all equal to each other i.e.  $d_4 = \dots = d_1 = d$ . We use the identity activation for the last layer (i.e.  $\phi_D(z) = z$ ) with the exception of the last experiment where we use a ReLU activation (i.e.  $\phi_D(z) = \max(0, z)$ ). We conducted our experiments in Python using the Tensorly library for the tensor decomposition in DeepTD [21].

	$d = 2, D = 12$			$d = 3, D = 8$		
	$N = 20$	$N = 50$	$N = 100$	$N = 20$	$N = 50$	$N = 100$
% Correct sign	0.83	0.95	0.93	0.65	0.87	0.94
Test loss (greedy)	0.53	0.45	0.43	0.62	0.48	0.42
Test loss (oracle)	0.50	0.44	0.41	0.58	0.44	0.40

Table 1: Test performance of DeepTD estimates (Test MSE =  $\mathbb{E}[(\mathbf{y} - \hat{f}_{\text{CNN}}(\mathbf{x}))^2] / \mathbb{E}[\mathbf{y}^2]$ ).

Greedy is the test loss when the signs of  $\hat{\mathbf{k}}^{(\ell)}$  are determined by Algorithm 1. Oracle is the test loss by picking the signs to ensure non-negative correlation with the ground truth kernels.

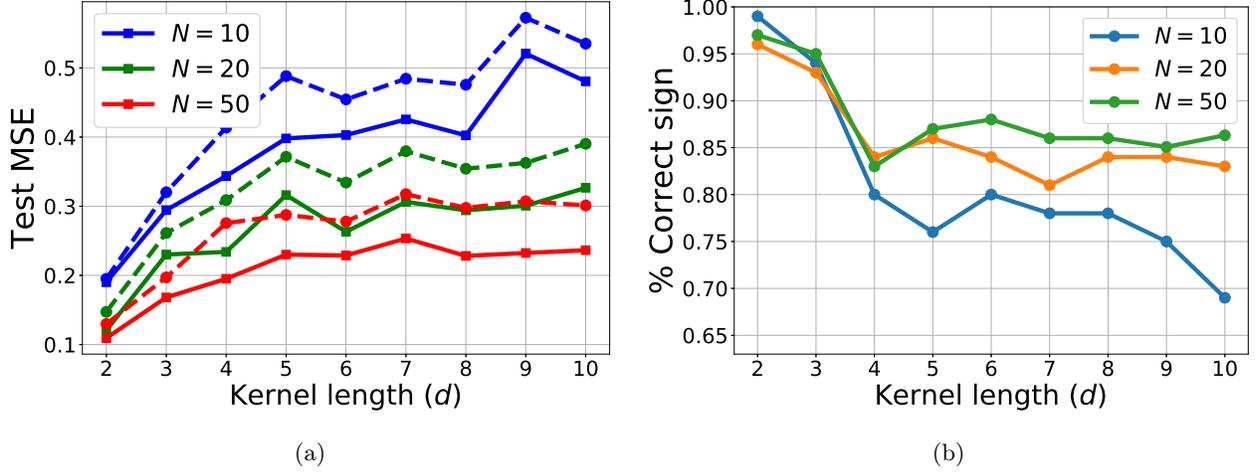


Figure 3: a) The solid lines are the Test MSE ( $\mathbb{E}[(\mathbf{y} - \hat{f}_{\text{CNN}}(\mathbf{x}))^2] / \mathbb{E}[\mathbf{y}^2]$ ) error when using the correct kernel signs (oracle). The dashed lines are the Test MSE by determining the kernel signs with Algorithm 1. b) Fraction of times Algorithm 1 successfully identified all signs.

Each curve in every figure is obtained by averaging 100 independent realizations of the same CNN learning procedure. Similar to our theory, we use Gaussian data points  $\mathbf{x}$  and ground truth labels  $\mathbf{y} = f_{\text{CNN}}(\mathbf{x})$ .

We conduct two sets of experiments: The first set focuses on larger values of depth  $D$  and the second set focuses on larger values of width  $d$ . In all experiments kernels are generated with random Gaussian entries and are normalized to have unit Euclidean norm. For the ReLU activation if one of the kernels have all negative entries, the output is trivially zero and learning is not feasible. To address this, we consider *operational* networks where at least 50% of the training labels are nonzero. Here, the number 50% is arbitrarily chosen and we verified that similar results hold for other values. To ensure the kernels do not have all negative entries and the network is operational we use a rejection sampling scheme. That is, if our generated network does not obey the 50% assumption we discard that experiment and generate a new one. In practice, this is mainly an issue for small  $d$  values only as the chance of an all-negative kernel is equal to  $2^{-d}$ . Finally, to study the effect of finite samples, we let the sample size grow proportional to the total degrees of freedom  $\sum_{\ell=1}^D d_\ell$ . In particular, we set an oversampling factor  $N = \frac{n}{\sum_{\ell=1}^D d_\ell}$  and carry out the experiments for  $N \in \{10, 20, 50, 100\}$ . While our theory requires  $N \gtrsim \log D$ , in our experiments, we typically observe that improvement is marginal after  $N = 50$ .

In Figure 2, we consider two networks with  $d = 2, D = 12$  and  $d = 3, D = 8$  configurations. We plot the absolute correlation between the ground truth and the estimates as a function of layer depth. For each hidden layer  $1 \leq \ell \leq D$ , our correlation measure ( $y$ -axis) is

$$\text{corr}(\hat{\mathbf{k}}^{(\ell)}, \mathbf{k}^{(\ell)}) = \left| \langle \hat{\mathbf{k}}^{(\ell)}, \mathbf{k}^{(\ell)} \rangle \right|.$$

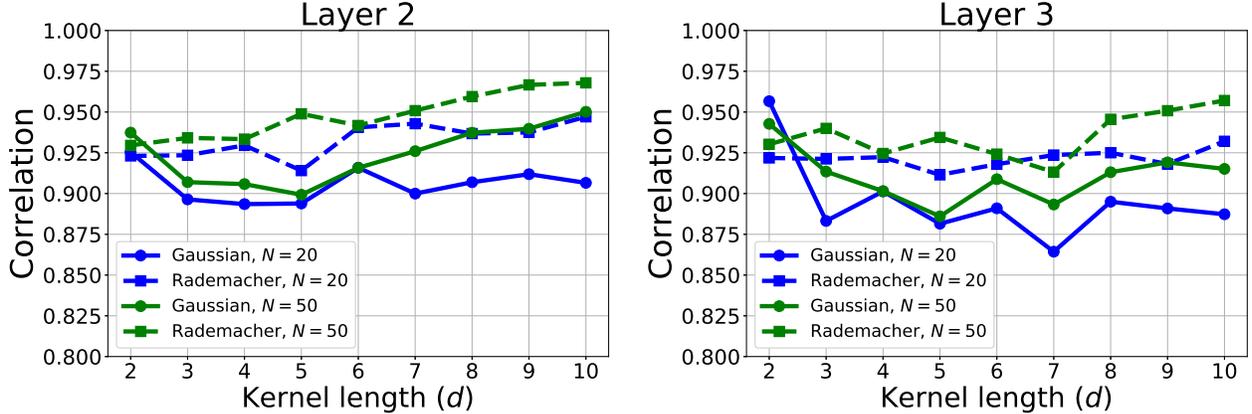


Figure 4: Correlations ( $\text{corr}(\hat{\mathbf{k}}^{(\ell)}, \mathbf{k}^{(\ell)}) = |\langle \hat{\mathbf{k}}^{(\ell)}, \mathbf{k}^{(\ell)} \rangle|$ ) when the ground truth kernels are Gaussian vs Rademacher. Gaussian is more spiky and results in lower correlation.

This number is between 0 and 1 as the kernels and their estimates both have unit norm. We observe that for both  $d = 2$  and  $d = 3$ , DeepTD consistently achieves correlation values above 75% for  $N = 20$ . While our theory requires  $d$  to scale quadratically with depth i.e.  $d \gtrsim D^2$ , we find that even small  $d$  values work well in our experiments. The effect of sample size becomes evident by comparing  $N = 20$  and  $N = 50$  for the input and output layers ( $\ell = 1, \ell = D$ ). In this case  $N = 50$  achieves perfect correlation. Interestingly, correlation values are smallest in the middle layers. In fact this even holds when  $N$  is large suggesting that the rank one approximation of the population tensor provides worst estimates for the middle layers.

Table 1 shows the test performance of DeepTD using an oracle and Algorithm 1. Our test loss is defined as the normalized regression loss

$$\text{Test MSE} = \frac{\mathbb{E} \left[ (\mathbf{y} - \hat{f}_{\text{CNN}}(\mathbf{x}))^2 \right]}{\mathbb{E}[\mathbf{y}^2]}, \quad (6.1)$$

where  $\hat{f}_{\text{CNN}}$  is our estimate of  $f_{\text{CNN}}(\cdot)$ . Here, the oracle knows the correct signs (i.e. guaranteed to have  $\langle \hat{\mathbf{k}}^{(\ell)}, \mathbf{k}^{(\ell)} \rangle \geq 0$  for all  $1 \leq \ell \leq D$ ) but the scale is found by solving a regression problem in a similar fashion to Line 12 of Algorithm 1. We observe that for  $N \geq 50$ , signs of all kernels are correctly identified with  $\geq 90\%$  accuracy. This table also demonstrates that there is a minor difference in test loss between the oracle and our proposed greedy algorithm. We would like to note that a test loss of  $\approx 40\%$  may appear large, especially as we achieve high kernel correlations  $\gtrsim 85\%$  (for  $N = 100$ ). We believe that this behavior is due to the depth of the network which aggregates the small estimation errors of individual layers and results in a larger error in the overall function estimate  $\hat{f}_{\text{CNN}}$ .

For the remaining experiments, we set  $D = 4$  and vary the width  $d$  from 2 to 10.  $N$  is varied between 10 to 50. Figure 3 summarizes the performance of the overall algorithm which is composed of DeepTD and Algorithm 1. In Figure 3a, we observe that small  $d$  values  $d = 2, 3$  achieve the best test performance. Performance appears to saturate after  $d \geq 7$  both for greedy algorithm (dashed lines) and oracle (solid lines) which has the correct sign information. The greedy algorithm is consistently worse than the oracle which indicates that it may have reached a local minima (this claim is verified on several examples where the oracle and greedy approaches have differing signs). While signs are often correctly identified in Figure 3b, the greedy algorithm fails in 15% of the instances even for large  $N$  which indicates the need for better heuristics to resolve sign and scale ambiguities.

Next we study our theoretical predictions regarding DeepTD. In particular, Figure 4 assesses the impact of diffuseness by plotting correlations for Layers 2 and 3. This is done by contrasting different kernel generation models, namely, kernels with i.i.d. standard normal entries and Rademacher ( $\pm 1$ ) entries (normalized to have

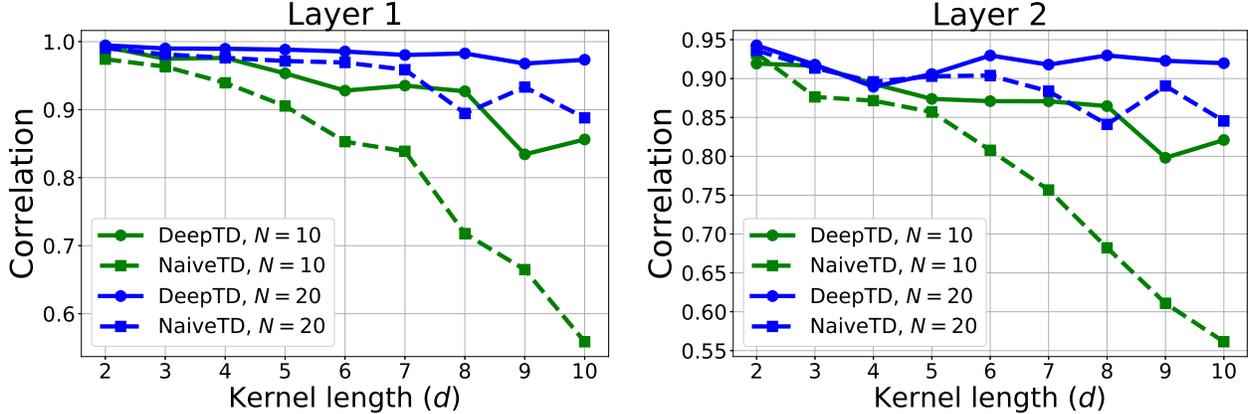


Figure 5: DeepTD estimate vs NaiveTD estimate when final activation is ReLU. Bias of NaiveTD results in significantly worse performance.

unit Euclidean norm). Rademacher kernels are strictly more diffused than Gaussian kernels as all entries have the same magnitude by construction. We observe that Rademacher kernels (dashed lines) achieve consistently higher correlations which is supported by our theory.

In Figure 5, we use a ReLU activation in the final layer and assess the impact of the centering procedure of the DeepTD algorithm which is a major theme throughout the paper. We define the NaiveTD algorithm which solves (3.2) without centering in the empirical tensor i.e.

$$\hat{\mathbf{k}}^{(1)}, \dots, \hat{\mathbf{k}}^{(D)} = \arg \max_{\mathbf{v}_1 \in \mathbb{R}^{d_1}, \mathbf{v}_2 \in \mathbb{R}^{d_2}, \dots, \mathbf{v}_D \in \mathbb{R}^{d_D}} \left\langle \frac{1}{n} \sum_{i=1}^n y_i \mathbf{X}_i, \bigotimes_{\ell=1}^D \mathbf{v}_\ell \right\rangle \quad \text{subject to} \quad \|\mathbf{v}_1\|_{\ell_2} = \dots = \|\mathbf{v}_D\|_{\ell_2} = 1. \quad (6.2)$$

Since the activation of the final layer is ReLU, the output has a clear positive bias in expectation which will help demonstrating the importance of centering. We find that for smaller oversampling factors of  $N = 10$  or  $N = 20$ , DeepTD has a visibly better performance compared with NaiveTD. The correlation difference is persistent among different layers (we plotted only Layers 1 and 2) and appears to grow with increase in the kernel size  $d$ .

Finally, in Figure 6, we assess the impact of activation nonlinearity by comparing the ReLU and identity activations in the final layer. We plot the first and final layer correlations for this setup. While the correlation performances of the first layer are essentially identical, the ReLU activation (dashed lines) achieves significantly lower correlation at the final layer. This is not surprising as the final layer passes through an additional nonlinearity.

## 7 Related work

Our work is closely related to the recent line of papers on neural networks as well as tensor decompositions. We briefly discuss this related literature.

**Neural networks:** Learning neural networks is a nontrivial task involving non-linearities and non-convexities. Consequently, existing theory works consider different algorithms, network structures and assumptions. A series of recent work focus on learning zero or one-hidden layer fully connected neural networks with random inputs and planted models [6, 13, 15, 17, 25, 27, 32, 33, 41]. Other recent publications [1, 11, 12, 27, 40] consider the problem of learning a convolutional neural network with 1-hidden layer. In particular, [1, 11, 12] focuses on learning non-overlapping networks as in this paper (albeit in the limit of infinite training data). The papers above either focus on characterizing the optimization landscape, or population landscape or providing exact convergence guarantees for gradient descent. In comparison, in this paper we focus on approximate convergence

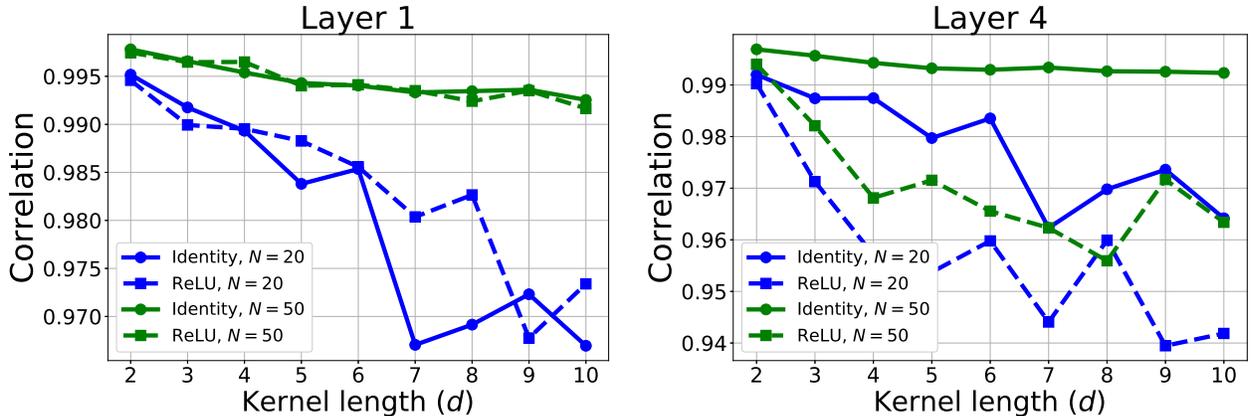


Figure 6: Comparison of performance of DeepTD when the final activation is ReLU in lieu of the identity activation.

guarantees using tensor decompositions for arbitrary deep networks. A few recent publications [30, 34, 35] consider the training problem when the network is over-parametrized and study the over-fitting ability of such networks. Closer to our work, [24] and [4] considers provable algorithms for deep networks using layer-wise algorithms. In comparison to our work, [24] applies to a very specific generative model that assumes a discrete data distribution and studies the population loss in lieu of the empirical loss. Arora et al. [4] studies deep models but uses activation functions that are not commonly used and assumes random network weights. In comparison, we work with practical activations and do not assume the weights are generated at random.

**Tensor decomposition:** Tensors are powerful tools to model a wide variety of big-data problems [2, 31]. Recent years have witnessed a growing interest in tensor decomposition techniques to extract useful latent information from the data [2, 14]. The connection between tensors and neural networks have been noticed by several papers [7, 8, 19, 20, 26]. Cohen et al. [7, 8] relate convolutional neural networks and tensor decompositions to provide insights on the expressivity of CNNs. Mondelli and Montanari [26] connects the hardness of learning shallow networks to tensor decompositions. Closer to this paper, Janzamin et al. [19] and Zhong et al. [41] apply tensor decomposition on *one*-hidden layer, fully connected networks to approximately learn the latent weight matrices.

## 8 Conclusion

In this paper we studied a multilayer CNN model with depth  $D$ . We assumed a non-overlapping structure where each layer has a single convolutional kernel and has stride length equal to the dimension of its kernel. We establish a connection between approximating the CNN kernels and higher order tensor decompositions. Based on this connection we proposed an algorithm for simultaneously learning all the kernels called the Deep Tensor Decomposition (DeepTD) algorithm. This algorithm builds a  $D$ -way tensor based on the training data and applies a rank one tensor factorization algorithm to this tensor to simultaneously estimate all the convolutional kernels. Assuming the input data is distributed i.i.d. according to a Gaussian model with corresponding output generated by a planted set of convolutional kernels we prove DeepTD can approximately learn all the kernels with a near minimal number of training data. A variety of numerical experiments complement our theoretical findings.

## 9 Proofs

In this section we will prove our main results. Throughout, for a random variable  $X$ , we use  $\mathbf{zm}(X)$  to denote  $X - \mathbb{E}[X]$ . Simply stated,  $\mathbf{zm}(X)$  is the *centered* version of  $X$ . For a random vector/matrix/tensor  $\mathbf{X}$ ,  $\mathbf{zm}(\mathbf{X})$  denotes the vector/matrix/tensor obtained by applying the  $\mathbf{zm}()$  operation to each entry. For a tensor  $\mathbf{T}$  we use  $\|\mathbf{T}\|_F$  to denote the square root of the sum of squares of the entries of the tensor. Stated differently, this is the Euclidean norm of a vector obtained by rearranging the entries of the tensor. Throughout we use  $c$ ,  $c_1$ ,  $c_2$ , and  $C$  to denote fixed numerical constants whose values may change from line to line. We begin with some useful definitions and lemmas.

### 9.1 Useful concentration lemmas and definitions

In this section we gather some useful definitions and well-known lemmas that will be used frequently throughout our concentration arguments.

**Definition 9.1 (Orlicz norms)** For a scalar random variable Orlicz- $a$  norm is defined as

$$\|X\|_{\psi_a} = \sup_{k \geq 1} k^{-1/a} (\mathbb{E}[|X|^k])^{1/k}$$

Orlicz- $a$  norm of a vector  $\mathbf{x} \in \mathbb{R}^p$  is defined as  $\|\mathbf{x}\|_{\psi_a} = \sup_{\mathbf{v} \in \mathcal{B}^p} \|\mathbf{v}^T \mathbf{x}\|_{\psi_a}$  where  $\mathcal{B}^p$  is the unit  $\ell_2$  ball. The sub-exponential norm is the function  $\|\cdot\|_{\psi_1}$  and the sub-gaussian norm the function  $\|\cdot\|_{\psi_2}$ .

We now state a few well-known results that we will use throughout the proofs. This results are standard and are stated for the sake of completeness. The first lemma states that the product of sub-gaussian random variables are sub-exponential.

**Lemma 9.2** Let  $X, Y$  be subgaussian random variables. Then  $\|XY\|_{\psi_1} \leq \|X\|_{\psi_2} \|Y\|_{\psi_2}$ .

The next lemma connects Orlicz norms of sum of random variables to the sum of the Orlicz norm of each random variable.

**Lemma 9.3** Suppose  $X, Y$  are random variables with bounded  $\|\cdot\|_{\psi_a}$  norm. Then  $\|X + Y\|_{\psi_a} \leq 2 \max\{\|X\|_{\psi_a}, \|Y\|_{\psi_a}\}$ . In particular  $\|X - \mathbb{E}X\|_{\psi_a} \leq 2\|X\|_{\psi_a}$ .

The lemma below can be easily obtained by combining the previous two lemmas.

**Lemma 9.4** Let  $X, Y$  be subgaussian random variables. Then  $\|\mathbf{zm}(XY)\|_{\psi_1} \leq 2\|X\|_{\psi_2} \|Y\|_{\psi_2}$ .

Finally, we need a few standard chaining definitions.

**Definition 9.5 (Admissible sequence [37])** Given a set  $T$  an admissible sequence is an increasing sequence  $\{A_n\}_{n=0}^{\infty}$  of partitions of  $T$  such that  $|A_n| \leq N_n$  where  $N_0 = 1$  and  $N_n = 2^{2^n}$  for all  $n \geq 1$ .

For the following discussion  $\Delta_d(A_n(t))$ , will be the diameter of the set  $S \in A_n$  that contains  $t$ , with respect to the  $d$  metric.

**Definition 9.6 ( $\gamma_a$  functional [37])** Given  $a > 0$ , and a metric space  $(T, d)$  we define

$$\gamma_a(T, d) = \inf \sup_{t \in T} \sum_{n \geq 0} 2^{n/a} \Delta_d(A_n(t)),$$

where the infimum is taken over all admissible sequences.

The following lemma upper bounds  $\gamma_\alpha$  functional with covering numbers of  $T$ . The reader is referred to Section 1.2 of [36], Equation (2.3) of [10], and Lemma D.17 of [27].

**Lemma 9.7 (Dudley's entropy integral)** Let  $N(\varepsilon)$  be the  $\varepsilon$  covering number of the set  $T$  with respect to the  $d$  metric. Then

$$\gamma_\alpha(T, d) \leq C_\alpha \int_0^\infty \log^{1/\alpha} N(\varepsilon) d\varepsilon,$$

where  $C_\alpha > 0$  depends only on  $\alpha > 0$ .

## 9.2 Concentration of the empirical tensor (Proof of Theorem 4.2)

To prove this theorem, first note that given labels  $\{y_i\}_{i=1}^n \sim y$  and their empirical average  $y_{\text{avg}} = n^{-1} \sum_{i=1}^n y_i$ , we have  $\mathbb{E}[y_{\text{avg}}] = \mathbb{E}[y]$ . Hence  $y - y_{\text{avg}} = \mathbf{z}\mathbf{m}(y - y_{\text{avg}})$  and we can rewrite the empirical tensor as follows

$$\begin{aligned} \mathbf{T}_n &:= \frac{1}{n} \sum_{i=1}^n (y_i - y_{\text{avg}}) \mathbf{X}_i \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(y_i - y_{\text{avg}}) \mathbf{X}_i \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(y_i) \mathbf{X}_i - \mathbf{z}\mathbf{m}(y_{\text{avg}}) \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \right). \end{aligned} \tag{9.1}$$

Recall that the population tensor is equal to  $\mathbf{T} := \mathbb{E}[y_i \mathbf{X}_i]$ . Furthermore,  $\mathbb{E}[\mathbb{E}[y_i] \mathbf{X}_i] = \mathbb{E}[y_i] \mathbb{E}[\mathbf{X}_i] = 0$ . Thus the population tensor can alternatively be written as  $\mathbf{T} = \mathbb{E}[\mathbf{z}\mathbf{m}(y_i) \mathbf{X}_i] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{z}\mathbf{m}(y_i) \mathbf{X}_i]$ . Combining the latter with (9.1) we conclude that

$$\begin{aligned} \mathbf{T}_n - \mathbf{T} &= \frac{1}{n} \sum_{i=1}^n (\mathbf{z}\mathbf{m}(y_i) \mathbf{X}_i - \mathbb{E}[\mathbf{z}\mathbf{m}(y_i) \mathbf{X}_i]) - \mathbf{z}\mathbf{m}(y_{\text{avg}}) \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \right), \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(\mathbf{z}\mathbf{m}(y_i) \mathbf{X}_i) - \mathbf{z}\mathbf{m}(y_{\text{avg}}) \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \right), \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(\mathbf{z}\mathbf{m}(y_i) \mathbf{X}_i) - \left( \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(y_i) \right) \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \right). \end{aligned}$$

Now using the triangular inequality for tensor spectral norm we conclude that

$$\|\mathbf{T}_n - \mathbf{T}\| \leq \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(\mathbf{z}\mathbf{m}(y_i) \mathbf{X}_i) \right\| + \left\| \left( \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(y_i) \right) \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \right) \right\|.$$

We now state two lemmas to bound each of these terms. The proofs of these lemmas are deferred to Sections 9.2.1 and 9.2.2.

**Lemma 9.8** For  $i = 1, 2, \dots, n$  let  $\mathbf{x}_i \in \mathbb{R}^p$  be i.i.d. random Gaussian vectors distributed as  $\mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ . Also let  $\mathbf{X}_i \in \mathbb{R}^{\otimes_{\ell=1}^D d_\ell}$  be the tensorized version of  $\mathbf{x}_i$  i.e.  $\mathbf{X}_i = \mathcal{T}(\mathbf{x}_i)$ . Finally, assume  $f : \mathbb{R}^p \mapsto \mathbb{R}$  is an  $L$  Lipschitz function. Furthermore, assume  $n \geq \left( \sum_{\ell=1}^D d_\ell \right) \log D$  and  $D \geq 2$ . Then

$$\mathbb{P} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(\mathbf{z}\mathbf{m}(f(\mathbf{x}_i)) \mathbf{X}_i) \right\| \leq \frac{c_1 L}{\sqrt{n}} \left( \sqrt{\left( \sum_{\ell=1}^D d_\ell \right) \log D} + t \right) \right\} \leq e^{-\min(t^2, t\sqrt{n})},$$

holds with  $c_1 > 0$  a fixed numerical constant.

**Lemma 9.9** Consider the setup of Lemma 9.8. Then

$$\mathbb{P} \left\{ \left\| \left( \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(f(\mathbf{x}_i)) \right) \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \right) \right\| \geq \frac{c_2 t_1 L}{n} \left( \sqrt{\left( \sum_{\ell=1}^D d_\ell \right) \log D} + t_2 \right) \right\} \leq 2 \left( e^{-t_1^2} + e^{-t_2^2} \right),$$

holds with  $c_2 > 0$  a fixed numerical constant.

Combining Lemma 9.8 with  $f = f_{\text{CNN}}()$ ,  $L = \prod_{\ell=1}^D \|\mathbf{k}^{(\ell)}\|_{\ell_2}$ , and  $c_1 = c/2$  together with Lemma 9.9 with  $t_1 = \sqrt{n}$ ,  $t_2 = t$ , and  $c_2 = c/2$  concludes the proof of Theorem 4.2. All that remains is to prove Lemmas 9.8 and 9.9 which are the subject of the next two sections.

### 9.2.1 Proof of Lemma 9.8

It is more convenient to carry out the steps of the proof on  $\sum_{i=1}^n \mathbf{zm}(\mathbf{zm}(f(\mathbf{x}_i))\mathbf{X}_i)$  in lieu of  $\frac{1}{n} \sum_{i=1}^n \mathbf{zm}(\mathbf{zm}(f(\mathbf{x}_i))\mathbf{X}_i)$ . The lemma trivially follows by a scaling by a factor  $1/n$ . We first write the tensor spectral norm as a supremum

$$\left\| \sum_{i=1}^n \mathbf{zm}(\mathbf{zm}(f(\mathbf{x}_i))\mathbf{X}_i) \right\| = \sup_{\mathbf{T} \in \mathcal{RO}} \left| \sum_{i=1}^n \langle \mathbf{zm}(\mathbf{zm}(f(\mathbf{x}_i))\mathbf{X}_i), \mathbf{T} \rangle \right|. \quad (9.2)$$

Let  $\mathbf{Y}_i = \mathbf{zm}(f(\mathbf{x}_i))\mathbf{X}_i$ . Define the random process  $g(\mathbf{T}) = \sum_{i=1}^n \langle \mathbf{zm}(\mathbf{Y}_i), \mathbf{T} \rangle$ . We claim that  $g(\mathbf{T})$  has a mixture of subgaussian and subexponential increments (see Definition 9.1 for subgaussian and subexponential random variables). Pick two tensors  $\mathbf{T}, \mathbf{H} \in \mathbb{R}^{\otimes_{\ell=1}^D d_\ell}$ . Increments of  $g$  satisfy the linear relation

$$g(\mathbf{T}) - g(\mathbf{H}) = \sum_{i=1}^n \langle \mathbf{zm}(\mathbf{Y}_i), \mathbf{T} - \mathbf{H} \rangle.$$

By construction  $\mathbb{E}[g(\mathbf{T}) - g(\mathbf{H})] = 0$ . We next claim that  $\mathbf{Y}_i$  is a sub-exponential vector. Consider a tensor  $\mathbf{T}$  with unit length  $\|\mathbf{T}\|_F = 1$  i.e. the sum of squares of entries are equal to one. We have  $\langle \mathbf{Y}_i, \mathbf{T} \rangle = \mathbf{zm}(f(\mathbf{x}_i)) \langle \mathbf{X}_i, \mathbf{T} \rangle$ .  $f(\mathbf{X}_i)$  is a Lipschitz function of a Gaussian random vector. Thus, by the concentration of Lipschitz functions of Gaussians we have

$$\mathbb{P}(|\mathbf{zm}(f(\mathbf{X}_i))| \geq t) \leq 2 \exp\left(-\frac{t^2}{2L^2}\right). \quad (9.3)$$

This immediately implies that  $\|\mathbf{zm}(f(\mathbf{X}_i))\|_{\psi_2} \leq cL$  for a fixed numerical constant  $c$ . Also note that  $\langle \mathbf{X}_i, \mathbf{T} \rangle \sim \mathcal{N}(0, 1)$  hence  $\|\langle \mathbf{X}_i, \mathbf{T} \rangle\|_{\psi_2} \leq c$ . These two identities combined with Lemma 9.4 implies a bound on the sub-exponential norm

$$\|\mathbf{zm}(\langle \mathbf{Y}_i, \mathbf{T} \rangle)\|_{\psi_1} \leq cL.$$

Next, we observe that  $g(\mathbf{T}) - g(\mathbf{H})$  is sum of  $n$  i.i.d. sub-exponentials each obeying  $\|\langle \mathbf{zm}(\mathbf{Y}_i), \mathbf{T} - \mathbf{H} \rangle\|_{\psi_1} \leq cL\|\mathbf{T} - \mathbf{H}\|_F$ . Applying a standard sub-exponential Bernstein inequality, we conclude that

$$\mathbb{P}(|g(\mathbf{T}) - g(\mathbf{H})| \geq t) \leq 2 \exp\left(-c \cdot \min\left(\frac{t}{L\|\mathbf{T} - \mathbf{H}\|_F}, \frac{t^2}{nL^2\|\mathbf{T} - \mathbf{H}\|_F^2}\right)\right), \quad (9.4)$$

holds with  $\gamma$  a fixed numerical constant. This tail bound implies that  $g$  is a mixed tail process that is studied by Talagrand and others [10, 37]. In particular, supremum of such processes are characterized in terms of a linear combination of Talagrand's  $\gamma_1$  and  $\gamma_2$  functionals (see Definition 9.6 as well as [36, 37] for an exposition). We pick the following distance metrics on tensors induced by the Frobenius norm:  $d_1(\mathbf{T}, \mathbf{H}) = L\|\mathbf{H} - \mathbf{T}\|_F/c$  and  $d_2(\mathbf{T}, \mathbf{H}) = \|\mathbf{H} - \mathbf{T}\|_F L\sqrt{n/c}$ . We can thus rewrite (9.4) in the form

$$\mathbb{P}\{|g(\mathbf{T}) - g(\mathbf{H})| \geq t\} \leq 2 \exp\left(-\min\left(\frac{t}{d_1(\mathbf{T}, \mathbf{H})}, \frac{t^2}{d_2^2(\mathbf{T}, \mathbf{H})}\right)\right),$$

which implies  $\mathbb{P}\{|g(\mathbf{T}) - g(\mathbf{H})| \geq \sqrt{t}d_2(\mathbf{T}, \mathbf{H}) + td_1(\mathbf{T}, \mathbf{H})\} \leq 2 \exp(-t)$ . Observe that the radius of  $\mathcal{RO}$  with respect to  $\|\cdot\|_F$  norm is 1 hence radius with respect to  $d_1, d_2$  metrics are  $L/c, L\sqrt{n/c}$  respectively. Applying Theorem 3.5 of Dirksen [10], we obtain

$$\mathbb{P}\left\{\sup_{\mathbf{T} \in \mathcal{RO}} |g(\mathbf{T})| \geq C\left(\gamma_2(\mathcal{RO}, d_2) + \gamma_1(\mathcal{RO}, d_1) + L\sqrt{un/c} + uL/c\right)\right\} \leq e^{-u}.$$

Observe that we can use the change of variable  $t = L \cdot \max(\sqrt{un}, u)$  to obtain

$$\mathbb{P}\left\{\sup_{\mathbf{T} \in \mathcal{RO}} |g(\mathbf{T})| \geq C\left(\gamma_2(\mathcal{RO}, d_2) + \gamma_1(\mathcal{RO}, d_1) + t\right)\right\} \leq \exp\left(-\min\left(\frac{t^2}{L^2n}, \frac{t}{L}\right)\right), \quad (9.5)$$

with some updated constant  $C > 0$ . To conclude, we need to bound the  $\gamma_2$  and  $\gamma_1$  terms. To achieve this we will upper bound the  $\gamma_\alpha$  functional in terms of Dudley's entropy integral which is stated in Lemma 9.7. First, let us find the  $\varepsilon$  covering number of  $\mathcal{RO}$ . Pick  $0 < \delta \leq 1$  coverings  $\mathcal{C}_\ell$  of the unit  $\ell_2$  balls  $\mathcal{B}^{d_\ell}$ . These covers have size at most  $(1 + 2/\delta)^{d_\ell}$ . Consider the set of rank 1 tensors  $\mathcal{C} = \mathcal{C}_1 \otimes \dots \otimes \mathcal{C}_D$  with size  $(1 + 2/\delta)^{\sum_{\ell=1}^D d_\ell}$ . For any  $\otimes_{\ell=1}^D \mathbf{v}_\ell \in \mathcal{RO}$ , we can pick  $\otimes_{\ell=1}^D \mathbf{u}_\ell \in \mathcal{C}$  satisfying  $\|\mathbf{v}_\ell - \mathbf{u}_\ell\|_{\ell_2} \leq \delta$  for all  $1 \leq \ell \leq D$ . This implies

$$\|(\{\mathbf{v}_\ell\}_{\ell=1}^D) - (\{\mathbf{u}_\ell\}_{\ell=1}^D)\|_F \leq \sum_{\ell=1}^D \|(\mathbf{v}_1, \dots, \mathbf{v}_\ell, \mathbf{u}_{\ell+1}, \dots, \mathbf{u}_D) - (\mathbf{v}_1, \dots, \mathbf{v}_{\ell-1}, \mathbf{u}_\ell, \dots, \mathbf{u}_D)\|_F \quad (9.6)$$

$$= \sum_{\ell=1}^D \|\mathbf{v}_\ell\|_{\ell_2} \dots \|\mathbf{v}_{\ell-1}\|_{\ell_2} \|\mathbf{v}_\ell - \mathbf{u}_\ell\|_{\ell_2} \|\mathbf{u}_{\ell+1}\|_{\ell_2} \dots \|\mathbf{u}_D\|_{\ell_2} \leq D\delta. \quad (9.7)$$

Denoting Frobenius norm covering number of  $\mathcal{RO}$  by  $N(\varepsilon)$ , this implies that, for  $0 < \varepsilon \leq 1$ ,

$$N(\varepsilon) \leq (1 + 2D/\varepsilon)^{\sum_{\ell=1}^D d_\ell}.$$

Clearly,  $N(\varepsilon) = 1$  for  $\varepsilon \geq 1$  by picking the cover  $\{0\}$ . Consequently,

$$\gamma_1(\mathcal{RO}, \|\cdot\|_F) \leq \int_0^1 \log N(\varepsilon) d\varepsilon \leq c \left( \sum_{\ell=1}^D d_\ell \right) \log D, \quad \gamma_2(\mathcal{RO}, \|\cdot\|_F) \leq \int_0^1 \sqrt{\log N(\varepsilon)} d\varepsilon \leq c \sqrt{\left( \sum_{\ell=1}^D d_\ell \right) \log D}. \quad (9.8)$$

Thus the metrics  $d_1, d_2$  metrics are  $\|\cdot\|_F$  norm scaled by a constant. Hence, their  $\gamma_\alpha$  functions are scaled versions of (9.8) given by

$$\gamma_1(\mathcal{RO}, d_1) \leq cL \left( \sum_{\ell=1}^D d_\ell \right) \log D, \quad \gamma_2(\mathcal{RO}, d_2) \leq cL \sqrt{n \left( \sum_{\ell=1}^D d_\ell \right) \log D}.$$

Now, observe that if  $n \geq \left( \sum_{\ell=1}^D d_\ell \right) \log D$ , we have  $\gamma_1(\mathcal{RO}, d_1) \leq cL \sqrt{n \left( \sum_{\ell=1}^D d_\ell \right) \log D}$ . Substituting these in (9.5) and using  $n \geq \left( \sum_{\ell=1}^D d_\ell \right) \log D$  we find

$$\mathbb{P} \left\{ \sup_{\mathbf{T} \in \mathcal{RO}} |g(\mathbf{T})| \geq C \left( L \sqrt{n \left( \sum_{\ell=1}^D d_\ell \right) \log D} + t \right) \right\} \leq e^{-\min\left(\frac{t^2}{L^2 n}, \frac{t}{L}\right)}.$$

Substituting  $t \rightarrow L\sqrt{nt}$  and recalling (9.2), concludes the proof.

### 9.2.2 Proof of Lemma 9.9

We first rewrite,

$$\left\| \left( \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(f(\mathbf{x}_i)) \right) \left( \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \right) \right\| = \left| \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(f(\mathbf{x}_i)) \right| \cdot \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \right\|. \quad (9.9)$$

As discussed in (9.3),  $\|\mathbf{z}\mathbf{m}(f(\mathbf{x}_i))\|_{\psi_2} \leq cL$  for a fixed numerical constant. Since  $\mathbf{x}_i$ 's are i.i.d the empirical sum  $\mathbf{f}_{\text{avg}} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(f(\mathbf{x}_i))$  obeys the bound  $\|\mathbf{f}_{\text{avg}}\|_{\psi_2} \leq cL/\sqrt{n}$  as well. Hence,

$$\mathbb{P} \left\{ |\mathbf{f}_{\text{avg}}| \geq c' \frac{t_1 L}{\sqrt{n}} \right\} \leq 2e^{-t_1^2}. \quad (9.10)$$

Also note that  $\frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{X}_i$  is a tensor with standard normal entries, applying [38, Theorem 1] we conclude that

$$\left\| \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{X}_i \right\| \leq c'' \left( \sqrt{\left( \sum_{\ell=1}^D d_\ell \right) \log D} + t_2 \right) \Rightarrow \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \right\| \leq c'' \frac{\sqrt{\left( \sum_{\ell=1}^D d_\ell \right) \log D} + t_2}{\sqrt{n}}, \quad (9.11)$$

holds with probability  $1 - 2e^{-t_2^2}$ . Combining (9.10) and (9.11) via the union bound together with (9.9) concludes the proof.

### 9.3 Rank one approximation of the population tensor (Proof of Theorem 4.7)

We begin the proof of this theorem by a few definitions regarding non-overlapping CNN models that simplify our exposition. For these definitions it is convenient to view non-overlapping CNNs as a tree with the root of the tree corresponding to the output of the CNN and the leaves corresponding to input features. In this visualization  $D - \ell$ th layer of the tree corresponds to the  $\ell$ th layer. Figure 7 depicts such a tree visualization along with the definitions discussed below.

**Definition 9.10 (Path vector)** A vector  $\mathbf{i} \in \mathbb{R}^{D+1}$  is a path vector if its zeroth coordinate satisfies  $1 \leq \mathbf{i}_0 \leq p$  and for all  $D - 1 \geq j \geq 0$ ,  $\mathbf{i}_{j+1}$  obeys  $\mathbf{i}_{j+1} = \lceil \mathbf{i}_j / d_j \rceil$ . This implies  $1 \leq \mathbf{i}_j \leq p_j$  and  $\mathbf{i}_D = 1$ . We note that in the tree visualization a path vector corresponds to a path connecting a leaf (input feature) to the root of the tree (output). We use  $\mathcal{I}$  to denote the set of all path vectors and note that  $|\mathcal{I}| = p$ . We also define  $\mathbf{i}(i) \in \mathcal{I}$  be the vector whose zeroth entry is  $\mathbf{i}_0 = i$ . Stated differently  $\mathbf{i}(i)$  is the path connecting the input feature  $i$  to the output. Given a path  $\mathbf{i}$  and a  $p$  dimensional vector  $\mathbf{v}$ , we define  $\mathbf{v}_{\mathbf{i}} := \mathbf{v}_{\mathbf{i}_0}$ . A sample path vector is depicted in bold in Figure 7 which corresponds to  $\mathbf{i} = (d_1, 1, 1, 1)$ .

**Definition 9.11 (Kernel and activation path gains)** Consider a CNN model of the form (2.1) with input  $\mathbf{x}$  and inputs of hidden units given by  $\{\bar{\mathbf{h}}^{(\ell)}\}_{\ell=1}^D$ . To any path vector  $\mathbf{i}$  connecting an input feature to the output we associate two path gains: a kernel path gain denoted by  $\mathbf{k}_{\mathbf{i}}$  and activation path gain denoted by  $\phi'_{\mathbf{i}}(\mathbf{x})$  defined as

$$\mathbf{k}_{\mathbf{i}} = \prod_{\ell=1}^D \mathbf{k}_{\text{mod}(\mathbf{i}_{\ell-1}, d_{\ell})}^{(\ell)} \quad \text{and} \quad \phi'_{\mathbf{i}}(\mathbf{x}) = \prod_{\ell=1}^D \phi'_{\ell}(\bar{\mathbf{h}}_{\mathbf{i}_{\ell}}^{(\ell)}),$$

where  $\text{mod}(a, b)$  denotes the remainder of dividing integer  $a$  by  $b$ . In words the kernel path gain is multiplication of the kernel weights along the path and the activation path gain is the multiplication of the derivatives of the activations evaluated at the hidden units along the path. For the path depicted in Figure 7 in bold the kernel path gain is equal  $\mathbf{k}_{\mathbf{i}} = \mathbf{k}_{d_1}^{(1)} \mathbf{k}_1^{(2)} \mathbf{k}_1^{(3)}$  and the activation path gain is equal to  $\phi'_{\mathbf{i}}(\mathbf{x}) = \phi'_1(\bar{h}_1^{(1)}) \phi'_2(\bar{h}_1^{(2)}) \phi'_3(\bar{h}_1^{(3)})$ .

**Definition 9.12 (CNN offsprings)** Consider a CNN model of the form (2.1) with input  $\mathbf{x}$  and inputs of hidden units given by  $\{\bar{\mathbf{h}}^{(\ell)}\}_{\ell=1}^D$ . We will associate a set  $\text{set}_{\ell}(i) \subset \{1, \dots, p\}$  to the  $i$ th hidden unit of layer  $\ell$  defined as  $\text{set}_{\ell}(i) = \{(i-1)r_{\ell} + 1, (i-1)r_{\ell} + 2, \dots, ir_{\ell}\}$  where  $r_{\ell} = p/p_{\ell}$ . By construction, this corresponds to the set of entries of the input data  $\mathbf{x}$  that  $\bar{\mathbf{h}}_i^{(\ell)}(\mathbf{x})$  is dependent on. In the tree analogy  $\text{set}_{\ell}(i)$  are the leaves of the tree connected to hidden unit  $i$  in the  $\ell$ th layer i.e. the set of offsprings of this hidden node. We depict  $\text{set}_2(p_2)$  which are the offsprings of the last hidden node in layer two in Figure 7.

We now will rewrite the population tensor in a form such that it is easier to see why it can be well approximated by a rank one tensor. Note that since the tensorization operation is linear the population tensor is equal to

$$\mathbf{T} = \mathbb{E}[f_{\text{CNN}}(\mathbf{x})\mathbf{X}] = \mathbb{E}[f_{\text{CNN}}(\mathbf{x})\mathcal{T}(\mathbf{x})] = \mathcal{T}(\mathbb{E}[f_{\text{CNN}}(\mathbf{x})\mathbf{x}]). \quad (9.12)$$

Define the vector  $\mathbf{g}^{\text{CNN}}$  to be the population gradient vector i.e.  $\mathbf{g}^{\text{CNN}} = \mathbb{E}[\nabla f_{\text{CNN}}(\mathbf{x})]$  and note that Stein's lemma combined with (9.12) implies that

$$\mathbf{T} = \mathcal{T}(\mathbb{E}[f_{\text{CNN}}(\mathbf{x})\mathbf{x}]) = \mathcal{T}(\mathbb{E}[\nabla f_{\text{CNN}}(\mathbf{x})]) = \mathcal{T}(\mathbf{g}^{\text{CNN}}). \quad (9.13)$$

Also note that

$$\frac{\partial f_{\text{CNN}}(\mathbf{x})}{\partial \mathbf{x}_i} = \mathbf{k}_{\mathbf{i}(i)} \phi'_{\mathbf{i}(i)}(\mathbf{x}).$$

Thus we have

$$\mathbf{g}_i^{\text{CNN}} = \mathbb{E}\left[\frac{\partial f_{\text{CNN}}(\mathbf{x})}{\partial \mathbf{x}_i}\right] = \mathbf{k}_{\mathbf{i}(i)} \mathbb{E}[\phi'_{\mathbf{i}(i)}(\mathbf{x})]. \quad (9.14)$$

$$\mathbf{x} = \mathbf{h}^{(0)} \in \mathbb{R}^p := \mathbb{R}^{p_0}$$

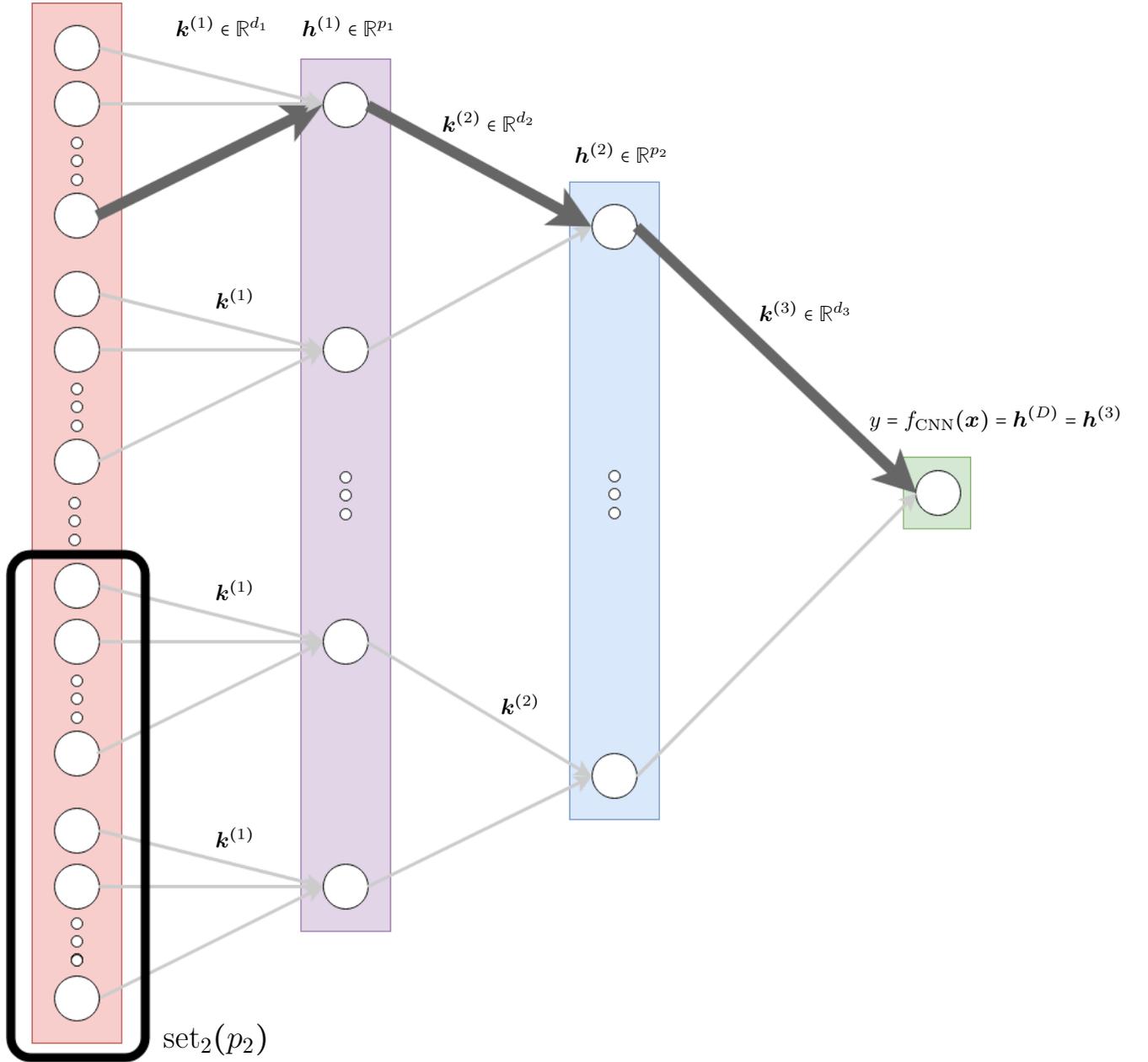


Figure 7: Tree visualization of a non-overlapping CNN model. The path in bold corresponds to the path vector  $\mathbf{i} = (d_1, 1, 1, 1)$  from Definition 9.10. For this path the kernel path gain is equal  $k_{\mathbf{i}} = \mathbf{k}_{d_1}^{(1)} \mathbf{k}_1^{(2)} \mathbf{k}_1^{(3)}$  and the activation path gain is equal to  $\phi'_{\mathbf{i}}(\mathbf{x}) = \phi'_1(\bar{h}_1^{(1)}) \phi'_2(\bar{h}_1^{(2)}) \phi'_3(\bar{h}_1^{(3)})$ . The set  $\text{set}_2(p_2)$  (offsprings of the last hidden node in layer two) is outlined.

Define a vector  $\mathbf{k} \in \mathbb{R}^p$  such that  $\mathbf{k}_i = \mathbf{k}_{\mathbf{i}(i)}$ . Since  $\mathbf{k}_{\mathbf{i}(i)}$  consists of the product of the kernel values across the path  $\mathbf{i}(i)$  it is easy to see that the tensor  $\mathbf{K} := \mathcal{T}(\mathbf{k})$  is equal to

$$\mathbf{K} = \bigotimes_{\ell=1}^D \mathbf{k}^{(\ell)}. \quad (9.15)$$

Similarly define the vector  $\mathbf{v} \in \mathbb{R}^p$  such that  $\mathbf{v}_i = E[\phi'_{\mathbf{i}(i)}(\mathbf{x})]$  and define the corresponding tensor  $\mathbf{V} = \mathcal{T}(\mathbf{v})$ . Therefore, (9.14) can be rewritten in the vector form  $\mathbf{g}^{\text{CNN}} = \mathbf{k} \odot \mathbf{v}$  where  $\mathbf{a} \odot \mathbf{b}$  denotes entry-wise (Hadamard) product between two vectors/matrices/tensors  $\mathbf{a}$  and  $\mathbf{b}$  of the same size. Thus using (9.13) the population tensor can alternatively be written as

$$\mathbf{T} = \mathcal{T}(\mathbf{g}^{\text{CNN}}) = \mathbf{K} \odot \mathbf{V} = \left( \bigotimes_{\ell=1}^D \mathbf{k}^{(\ell)} \right) \odot \mathbf{V}.$$

Therefore, the population tensor  $\mathbf{T}$  is the outer product of the convolutional kernels whose entries are masked with another tensor  $\mathbf{V}$ . If the entries of the tensor  $\mathbf{V}$  were all the same the population tensor would be exactly rank one with the factors revealing the convolutional kernel. One natural choice for approximating the population tensor with a rank one matrix is to replace the masking tensor  $\mathbf{V}$  with a scalar. That is, use the approximation  $\mathbf{T} \approx c \bigotimes_{\ell=1}^D \mathbf{k}^{(\ell)}$ . Recall that  $\mathbf{L}_{\text{CNN}} := \alpha_{\text{CNN}} \bigotimes_{\ell=1}^D \mathbf{k}^{(\ell)}$  is exactly such an approximation with  $c$  set to  $\alpha_{\text{CNN}}$  given by

$$\alpha_{\text{CNN}} = \prod_{\ell=1}^D \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})} [\phi'_{\ell}(\bar{\mathbf{h}}_{\mathbf{i}_{\ell}}^{(\ell)})].$$

To characterize the quality of this approximation note that

$$\begin{aligned} \|\mathbf{T} - \mathbf{L}_{\text{CNN}}\| &\stackrel{(a)}{\leq} \|\mathbf{T} - \mathbf{L}_{\text{CNN}}\|_F, \\ &\stackrel{(b)}{=} \|\mathbf{K} \odot (\mathbf{V} - \alpha_{\text{CNN}})\|_F, \\ &\stackrel{(c)}{=} \|\mathbf{k} \odot (\mathbf{v} - \alpha_{\text{CNN}})\|_{\ell_2}, \\ &\stackrel{(d)}{\leq} \|\mathbf{k}\|_{\ell_2} \|\mathbf{v} - \alpha_{\text{CNN}}\|_{\ell_{\infty}}, \\ &\stackrel{(e)}{=} \prod_{\ell=1}^D \|\mathbf{k}^{(\ell)}\|_{\ell_2} \|\mathbf{v} - \alpha_{\text{CNN}}\|_{\ell_{\infty}}. \end{aligned}$$

Here, (a) follows from the fact for a tensor, its spectral norm is smaller than its Frobenius norm, (b) from the definitions of  $\mathbf{T}$  and  $\mathbf{L}_{\text{CNN}}$ , (c) from the fact that  $\mathbf{K} = \mathcal{T}(\mathbf{k})$  and  $\mathbf{V} = \mathcal{T}(\mathbf{v})$ , (d) from the fact that  $\|\mathbf{a} \odot \mathbf{v}\|_{\ell_2} \leq \|\mathbf{a}\|_{\ell_2} \|\mathbf{v}\|_{\ell_{\infty}}$ , and (e) from the fact that the Euclidean norm of the kronecker product of vectors is equal to the product of the Euclidean norm of the individual vectors. As a result of the latter inequality to prove Theorem 4.7 it suffices to show that

$$\|\mathbf{v} - \alpha_{\text{CNN}}\|_{\ell_{\infty}} \leq \sqrt{8\pi} \mu S \cdot \sup_{\ell} \prod_{i=1}^{\ell} \|\mathbf{k}^{(i)}\|_{\ell_2} \cdot \frac{D}{\sqrt{\min_{\ell} d_{\ell}}}. \quad (9.16)$$

In the next lemma we prove a stronger statement.

**Lemma 9.13** *Assume the activations are  $S$ -smooth. Also consider a vector  $\mathbf{v} \in \mathbb{R}^p$  with entries  $\mathbf{v}_i = E_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})}[\phi'_{\mathbf{i}(i)}(\mathbf{x})]$  and  $\alpha_{\text{CNN}} = \prod_{\ell=1}^D \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})}[\phi'_{\ell}(\bar{\mathbf{h}}_{\mathbf{i}_{\ell}}^{(\ell)})]$  we have*

$$|\mathbf{v}_i - \alpha_{\text{CNN}}| \leq \kappa_i := \sqrt{8\pi} S \sum_{\ell=1}^D \left| \mathbf{k}_{\text{mod}(\mathbf{i}_{\ell-1}, d_{\ell})}^{(\ell)} \right| \prod_{i=1}^{\ell-1} \|\mathbf{k}^{(i)}\|_{\ell_2}.$$

Here,  $\mathbf{i}$  is the vector path that starts at input feature  $i$ .

Before proving this lemma let us explain how (9.16) follows from this lemma. To show this we use the kernel diffuseness assumption introduced in Definition 4.6. This definition implies that  $\left| \mathbf{k}_{\text{mod}(\mathbf{i}_{\ell-1}, d_\ell)}^{(\ell)} \right| \leq \|\mathbf{k}^{(\ell)}\|_{\ell_\infty} \leq \frac{\mu}{\sqrt{d_\ell}} \|\mathbf{k}^{(\ell)}\|_{\ell_2}$ . Thus we have

$$\begin{aligned}
\kappa_i &:= \sqrt{8\pi} S \sum_{\ell=1}^D |\mathbf{k}_{\text{mod}(\mathbf{i}_{\ell-1}, d_\ell)}^{(\ell)}| \prod_{i=1}^{\ell-1} \|\mathbf{k}^{(i)}\|_{\ell_2}, \\
&\leq \sqrt{8\pi} \mu S \sum_{\ell=1}^D \frac{1}{\sqrt{d_\ell}} \|\mathbf{k}^{(\ell)}\|_{\ell_2} \prod_{i=1}^{\ell-1} \|\mathbf{k}^{(i)}\|_{\ell_2}, \\
&= \sqrt{8\pi} \mu S \sum_{\ell=1}^D \frac{1}{\sqrt{d_\ell}} \prod_{i=1}^{\ell} \|\mathbf{k}^{(i)}\|_{\ell_2}, \\
&\leq \sqrt{8\pi} \mu D S \cdot \sup_{\ell} \frac{\prod_{i=1}^{\ell} \|\mathbf{k}^{(i)}\|_{\ell_2}}{\sqrt{d_\ell}}, \\
&\leq \sqrt{8\pi} \mu S \cdot \sup_{\ell} \prod_{i=1}^{\ell} \|\mathbf{k}^{(i)}\|_{\ell_2} \cdot \frac{D}{\sqrt{\min_{\ell} d_\ell}}.
\end{aligned}$$

This completes the proof of Theorem 4.7. All that remains is to prove Lemma 9.13 which is the subject of the next section.

### 9.3.1 Proof of Lemma 9.13

To bound the difference between  $\mathbf{v}_i$  and  $\alpha_{\text{CNN}}$  consider the path  $\mathbf{i} = \mathbf{i}(i)$  and define the variables  $\{a_i\}_{i=0}^D$  as

$$a_i = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})} \left[ \prod_{\ell=1}^i \phi'_\ell(\bar{\mathbf{h}}_{\mathbf{i}_\ell}^{(\ell)}) \right] \prod_{\ell=i+1}^D \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})} [\phi'_\ell(\bar{\mathbf{h}}_{\mathbf{i}_\ell}^{(\ell)})].$$

Note that  $a_D = \mathbf{v}_i$  and  $a_0 = \alpha_{\text{CNN}}$ . To bound the difference  $\mathbf{v}_i - \alpha_{\text{CNN}} = a_D - a_0$  we use a telescopic sum

$$|a_D - a_0| \leq \sum_{\ell=0}^{D-1} |a_{\ell+1} - a_\ell|. \quad (9.17)$$

We thus focus on bounding each of the summands  $|a_\ell - a_{\ell-1}|$ . Setting  $\gamma_\ell = \prod_{i=\ell}^D \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})} [\phi'_i(\bar{\mathbf{h}}_{\mathbf{i}_i}^{(i)})]$ , this can be written as

$$a_\ell - a_{\ell-1} = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})} \left[ (\phi'_\ell(\bar{\mathbf{h}}_{\mathbf{i}_\ell}^{(\ell)}) - \mathbb{E}[\phi'_\ell(\bar{\mathbf{h}}_{\mathbf{i}_\ell}^{(\ell)})]) \prod_{i=1}^{\ell-1} \phi'_i(\bar{\mathbf{h}}_{\mathbf{i}_i}^{(i)}) \right] \gamma_{\ell+1}.$$

Using  $|\gamma_\ell| \leq 1$  (which follows from the assumption that the activations are 1-Lipschitz), it suffices to bound

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})} \left[ (\phi'_\ell(\bar{\mathbf{h}}_{\mathbf{i}_\ell}^{(\ell)}) - \mathbb{E}[\phi'_\ell(\bar{\mathbf{h}}_{\mathbf{i}_\ell}^{(\ell)})]) \theta_{\ell-1} \right] \quad \text{where} \quad \theta_\ell = \prod_{i=1}^{\ell} \phi'_i(\bar{\mathbf{h}}_{\mathbf{i}_i}^{(i)}). \quad (9.18)$$

To this aim we state two useful lemmas whose proofs are deferred to Sections 9.3.1.1 and 9.3.1.2.

**Lemma 9.14** *Let  $X, Y, Z$  be random variables where  $X$  is independent of  $Z$ . Let  $f$  be an  $L$ -Lipschitz function. Then*

$$|\mathbb{E}[f(X+Y)Z] - \mathbb{E}[f(X+Y)]\mathbb{E}[Z]| \leq L \mathbb{E}[|\mathbf{z}\mathbf{m}(Y)|(|Z| + |\mathbb{E}[Z]|)]. \quad (9.19)$$

Furthermore if  $|Z| \leq 1$ , then

$$|\mathbb{E}[f(X+Y)Z] - \mathbb{E}[f(X+Y)]\mathbb{E}[Z]| \leq 2L \mathbb{E}[|\mathbf{z}\mathbf{m}(Y)|]. \quad (9.20)$$

**Lemma 9.15**  $\bar{\mathbf{h}}_i^{(\ell)}(\mathbf{x})$  (and  $\mathbf{h}_{i_\ell}^{(\ell)}(\mathbf{x})$ ) is a deterministic function of the entries of  $\mathbf{x}$  indexed by  $\text{set}_\ell(i)$ . In other words, there exists a function  $f$  such that  $\bar{\mathbf{h}}_i^{(\ell)}(\mathbf{x}) = f(\mathbf{x}_{\text{set}_\ell(i)})$ .

With these lemmas in-hand we return to bounding (9.18). To this aim we decompose  $\bar{\mathbf{h}}_{\mathbf{i}_\ell}^{(\ell)}$  as follows

$$\bar{\mathbf{h}}_{\mathbf{i}_\ell}^{(\ell)} = \sum_{i=1}^{d_\ell} \mathbf{k}_i^{(\ell)} \mathbf{h}_{d_\ell(\mathbf{i}_{\ell-1})+i}^{(\ell-1)} = \mathbf{k}_{\text{mod}(\mathbf{i}_{\ell-1}, d_\ell)}^{(\ell)} \mathbf{h}_{\mathbf{i}_{\ell-1}}^{(\ell-1)} + \mathbf{r},$$

where the  $\mathbf{r}$  term is the contribution of the entries of  $\mathbf{h}^{(\ell-1)}$  other than  $\mathbf{i}_{\ell-1}$ . By the non-overlapping assumption,  $\mathbf{r}$  is independent of  $\theta_{\ell-1}$  as well as  $\mathbf{h}_{\mathbf{i}_{\ell-1}}^{(\ell-1)}$  (see Lemma 9.15). In particular,  $\mathbf{h}_{\mathbf{i}_{\ell-1}}^{(\ell-1)}$  and  $\theta_{\ell-1}$  is a function of  $\mathbf{x}_{\text{set}_{\ell-1}(\mathbf{i}_{\ell-1})}$  whereas  $\mathbf{r}$  is a function of the entries over the complement  $\text{set}_\ell(\mathbf{i}_\ell) - \text{set}_{\ell-1}(\mathbf{i}_{\ell-1})$ . With these observations, applying Lemma 9.14 with  $f = \phi'_\ell$ ,  $X = \mathbf{r}$ ,  $Y = \mathbf{k}_{\text{mod}(\mathbf{i}_{\ell-1}, d_\ell)}^{(\ell)} \mathbf{h}_{\mathbf{i}_{\ell-1}}^{(\ell-1)}$ ,  $Z = \theta_{\ell-1}$  and using the fact that  $|\theta_{\ell-1}| \leq 1$  which holds due to 1-Lipschitzness of  $\sigma_\ell$ 's, we conclude that

$$|\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(0, \mathbf{I})}[(\phi'_\ell(\bar{\mathbf{h}}_{\mathbf{i}_\ell}^{(\ell)}) - \mathbb{E}[\phi'_\ell(\bar{\mathbf{h}}_{\mathbf{i}_\ell}^{(\ell)})])\theta_{\ell-1}]| \leq 2S \mathbb{E}|\mathbf{z}\mathbf{m}(Y)|.$$

Here,  $S$  is the smoothness of  $\sigma_\ell$  and Lipschitz constant of  $\phi'_\ell$ . To conclude, we need to assess the  $\mathbb{E}|\mathbf{z}\mathbf{m}(Y)|$  term. Now note that starting from  $\mathbf{x}$ , each entry of  $\mathbf{h}^{(\ell-1)}$  is obtained by applying a sequence of inner products with  $\{\mathbf{k}^{(i)}\}_{i=1}^{\ell-1}$  and activations  $\sigma_\ell(\cdot)$ , which implies  $\mathbf{h}_{\mathbf{i}_{\ell-1}}^{(\ell-1)}$  is a  $\prod_{i=1}^{\ell-1} \|\mathbf{k}^{(i)}\|_{\ell_2}$ -Lipschitz function of  $\text{set}_{\ell-1}(\mathbf{i}_{\ell-1})$ . This implies  $Y$  is a Lipschitz function of a Gaussian vector with Lipschitz constant  $L_Y = |\mathbf{k}_{\text{mod}(\mathbf{i}_{\ell-1}, d_\ell)}^{(\ell)}| \prod_{i=1}^{\ell-1} \|\mathbf{k}^{(i)}\|_{\ell_2}$ . Hence,  $\mathbf{z}\mathbf{m}(Y)$  obeys the tail bound

$$\mathbb{P}(|\mathbf{z}\mathbf{m}(Y)| \geq t) \leq 2 \exp\left(-\frac{t^2}{2L_Y^2}\right).$$

Using a standard integration by parts argument the latter implies that

$$\mathbb{E}|\mathbf{z}\mathbf{m}(Y)| \leq \sqrt{2\pi}L_Y.$$

Thus,

$$|a_\ell - a_{\ell-1}| \leq |\mathbf{k}_{\text{mod}(\mathbf{i}_{\ell-1}, d_\ell)}^{(\ell)}| \prod_{i=1}^{\ell-1} \|\mathbf{k}^{(i)}\|_{\ell_2}.$$

concluding the upper-bound on each summand of (9.17). Combining such upper bounds (9.17) implies

$$|a_D - a_0| \leq \sqrt{8\pi}S \sum_{\ell=1}^D |\mathbf{k}_{\text{mod}(\mathbf{i}_{\ell-1}, d_\ell)}^{(\ell)}| \prod_{i=1}^{\ell-1} \|\mathbf{k}^{(i)}\|_{\ell_2} := \kappa_{\mathbf{i}}.$$

This concludes the proof of Lemma 9.13.

### 9.3.1.1 Proof of Lemma 9.14

Using the independence of  $X, Z$ , we can write

$$\begin{aligned} \mathbb{E}[f(X+Y)Z] &= \mathbb{E}[f(X + \mathbb{E}[Y])Z] + \mathbb{E}[(f(X+Y) - f(X + \mathbb{E}[Y]))Z] \\ &= \mathbb{E}[f(X + \mathbb{E}[Y])] \mathbb{E}[Z] + \mathbb{E}[(f(X+Y) - f(X + \mathbb{E}[Y]))Z] \\ &= \mathbb{E}[f(X+Y)] \mathbb{E}[Z] + \mathbb{E}[f(X + \mathbb{E}[Y]) - f(X+Y)] \mathbb{E}[Z] + \mathbb{E}[(f(X+Y) - f(X + \mathbb{E}[Y]))Z]. \end{aligned}$$

This implies

$$\mathbb{E}[f(X+Y)Z] - \mathbb{E}[f(X+Y)] \mathbb{E}[Z] = \mathbb{E}[f(X + \mathbb{E}[Y]) - f(X+Y)] \mathbb{E}[Z] + \mathbb{E}[(f(X+Y) - f(X + \mathbb{E}[Y]))Z].$$

Now, using Lipschitzness of  $f$ , we deterministically have that  $|f(X + \mathbb{E}[Y]) - f(X + Y)| \leq L|Y - \mathbb{E}[Y]| = L|\mathbf{z}\mathbf{m}(Y)|$ . Similarly,  $|(f(X + Y) - f(X + \mathbb{E}[Y]))Z| \leq L|\mathbf{z}\mathbf{m}(Y)Z|$ . Taking absolute values we arrive at

$$|\mathbb{E}[f(X + \mathbb{E}[Y]) - f(X + Y)]\mathbb{E}[Z] + \mathbb{E}[(f(X + Y) - f(X + \mathbb{E}[Y]))Z]| \leq L\mathbb{E}[|\mathbf{z}\mathbf{m}(Y)|(|Z| + |\mathbb{E}[Z]|)].$$

This immediately implies (9.19). If  $|Z| \leq 1$  almost surely, we have  $\mathbb{E}[|\mathbf{z}\mathbf{m}(Y)Z|] \leq \mathbb{E}[|\mathbf{z}\mathbf{m}(Y)|]$  and  $|\mathbb{E}[Z]| \leq 1$  which yields the  $2L\mathbb{E}[|\mathbf{z}\mathbf{m}(Y)|]$  upper bound.

### 9.3.1.2 Proof of Lemma 9.15

Informally, this lemma is obvious via the tree visualization. To formally prove this lemma we use an induction argument. For  $\bar{\mathbf{h}}^{(1)}$  the result is trivial because  $\bar{\mathbf{h}}_i^{(1)} = \langle \mathbf{k}^{(1)}, \mathbf{x}(i) \rangle$  which is a weighted sum of entries corresponding to  $\text{set}_1(i)$ . Suppose the claim holds for all layers less than or equal to  $\ell - 1$  and  $\bar{\mathbf{h}}_i^{(\ell-1)} = f_{\ell-1}(\mathbf{x}_{\text{set}_{\ell-1}(i)})$ . For layer  $\ell$ , we can use the fact that  $\text{set}_\ell(i) = \bigcup_{j=1}^{d_\ell} \text{set}_{\ell-1}((i-1)d_\ell + j)$  to conclude that

$$\begin{aligned} \bar{\mathbf{h}}_i^{(\ell)} &= \sum_{j=1}^{d_\ell} \mathbf{k}_j^{(\ell)} \sigma_{\ell-1}(\bar{\mathbf{h}}_{(i-1)d_\ell + j}^{(\ell-1)}) \\ &= \sum_{j=1}^{d_\ell} \mathbf{k}_j^{(\ell)} \sigma_{\ell-1}(f_{\ell-1}(\mathbf{x}_{\text{set}_{\ell-1}((i-1)d_\ell + j)})) := f_\ell(\mathbf{x}_{\text{set}_\ell(i)}). \end{aligned}$$

The latter is clearly a deterministic function of  $\mathbf{x}_{\text{set}_\ell(i)}$ . Also it is independent of entry  $i$  because it simply chunks the vector  $\mathbf{x}_{\text{set}_\ell(i)}$  into  $d_\ell$  sub-vectors and returns a sum of weighted functions of these sub-vectors. Here, the weights are the entries of  $\mathbf{k}^{(\ell)}$  and the functions are given by  $\sigma_{\ell-1}(f_{\ell-1}(\cdot))$  (also note that the activation output is simply  $\mathbf{h}_i^{(\ell)} = \sigma_\ell(f_\ell(\mathbf{x}_{\text{set}_\ell(i)}))$ ).

## 9.4 Proofs for learning the convolutional kernels (Proof of Theorem 4.8)

The first part of the theorem follows trivially by combining Theorems 4.2 and 4.7. To translate a bound on the tensor spectral norm of  $\mathbf{T}_n - \mathbf{L}_{\text{CNN}}$  to a bound on learning the kernels, requires a perturbation argument for tensor decompositions. This is the subject of the next lemma.

**Lemma 9.16** *Let  $\mathbf{L} = \gamma \bigotimes_{\ell=1}^D \mathbf{v}_\ell$  be a rank one tensor with  $\{\mathbf{v}_i\}_{i=1}^D$  vectors of unit norm. Also assume  $\mathbf{E}$  is a perturbation tensor obeying  $\|\mathbf{E}\| \leq \delta$ . Set*

$$\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_D = \arg \max_{\bar{\mathbf{u}}_1, \bar{\mathbf{u}}_2, \dots, \bar{\mathbf{u}}_D} \left\langle \mathbf{L}, \bigotimes_{\ell=1}^D \bar{\mathbf{u}}_\ell \right\rangle \quad \text{subject to} \quad \|\bar{\mathbf{u}}_1\|_{\ell_2} = \|\bar{\mathbf{u}}_2\|_{\ell_2} = \dots = \|\bar{\mathbf{u}}_D\|_{\ell_2} = 1. \quad (9.21)$$

Then we have

$$\prod_{i=1}^D |\mathbf{u}_i^* \mathbf{v}_i| \geq 1 - 2\delta/\gamma.$$

The proof of Theorem 4.8 is complete by applying Lemma 9.16 above with  $\mathbf{v}_\ell = \mathbf{k}^{(\ell)}$ ,  $\mathbf{u}_\ell = \hat{\mathbf{k}}^{(\ell)}$ ,  $\gamma = \alpha_{\text{CNN}}$  and  $\mathbf{E} = \mathbf{T}_n - \mathbf{L}_{\text{CNN}}$ . All that remains is to prove Lemma 9.16 which is the subject of the next section.

### 9.4.1 Proof of Lemma 9.16

To prove this lemma first note that for any two rank one tensors we have

$$\left\langle \bigotimes_{\ell=1}^D \mathbf{u}_\ell, \bigotimes_{\ell=1}^D \mathbf{v}_\ell \right\rangle = \prod_{i=1}^D \langle \mathbf{u}_i, \mathbf{v}_i \rangle.$$

Using this equality together with the fact that the vectors  $\{\mathbf{u}_\ell\}_{\ell=1}^D$  are a maximizer for (9.21) we conclude that

$$\begin{aligned} \left\langle \mathbf{L} + \mathbf{E}, \bigotimes_{\ell=1}^D \mathbf{v}_\ell \right\rangle &\leq \left\langle \mathbf{L} + \mathbf{E}, \bigotimes_{\ell=1}^D \mathbf{u}_\ell \right\rangle, \\ &\leq \left| \left\langle \mathbf{L}, \bigotimes_{\ell=1}^D \mathbf{u}_\ell \right\rangle \right| + \left| \left\langle \mathbf{E}, \bigotimes_{\ell=1}^D \mathbf{u}_\ell \right\rangle \right|, \\ &\leq \gamma \prod_{i=1}^D |\langle \mathbf{u}_i, \mathbf{v}_i \rangle| + \delta. \end{aligned} \tag{9.22}$$

Furthermore, note that

$$\left\langle \mathbf{L} + \mathbf{E}, \bigotimes_{\ell=1}^D \mathbf{v}_\ell \right\rangle = \gamma + \left\langle \mathbf{E}, \bigotimes_{\ell=1}^D \mathbf{v}_\ell \right\rangle \geq \gamma - \delta. \tag{9.23}$$

Combining (9.22) and (9.23) we conclude that

$$\gamma \prod_{i=1}^D |\langle \mathbf{u}_i, \mathbf{v}_i \rangle| + \delta \geq \gamma - \delta,$$

concluding the proof.

## 9.5 Generalization bounds for CERM

In this section we prove a generalization bound for Centered Empirical Risk Minimization (CERM) (5.1). The following theorem shows that using a finite sample size  $n$ , CERM is guaranteed to choose a function close to population's minimizer. For the sake of this section  $\|f\|_{L_\infty}$  will be the Lipschitz constant of a function.

**Theorem 9.17** *Let  $\mathcal{F}$  be a class of functions  $f : \mathbb{R}^p \rightarrow \mathbb{R}$ . Suppose  $\sup_{f \in \mathcal{F}} \|f\|_{L_\infty} \leq R$ . Let  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim (\mathbf{x}, y)$  be i.i.d. data points where  $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_p)$  and  $y$  is so that  $y - \mathbb{E}[y]$  is  $K$  subgaussian. Suppose  $\mathcal{F}$  has  $\|\cdot\|_{L_\infty}$ ,  $\delta$ -covering bound obeying  $\log N_\delta \leq s \log \frac{C}{\delta}$  for some constants  $s \geq 1, C \geq R$ . Given  $\varepsilon \leq \bar{K} = K + R$ , suppose  $n \geq \mathcal{O}(\max\{\varepsilon^{-1}, \varepsilon^{-2}\} \bar{K}^2 s \log \frac{C' p \bar{K}}{\varepsilon})$  for some  $C' > 0$ . Then the CERM output (5.1) obeys*

$$\mathbb{E}[(\hat{f}(\mathbf{x}) - y - \mathbb{E}[(\hat{f}(\mathbf{x}) - y)])^2] \leq \min_{f \in \mathcal{F}} \mathbb{E}[(f(\mathbf{x}) - y - \mathbb{E}[(f(\mathbf{x}) - y)])^2] + \varepsilon. \tag{9.24}$$

with probability  $1 - \exp(-\mathcal{O}(n)) - 4n \exp(-p)$ .

**Proof** Consider the centered empirical loss that can alternatively be written in the form

$$\begin{aligned} E(f) &= \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(y_i - f(\mathbf{x}_i))^2 - \frac{1}{n^2} \left( \sum_{i=1}^n \mathbf{z}\mathbf{m}(y_i - f(\mathbf{x}_i)) \right)^2 - \mathbb{E}[\mathbf{z}\mathbf{m}(f(\mathbf{x}) - y)^2]. \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(\mathbf{z}\mathbf{m}(y_i - f(\mathbf{x}_i))^2) - \frac{1}{n^2} \left( \sum_{i=1}^n \mathbf{z}\mathbf{m}(y_i - f(\mathbf{x}_i)) \right)^2 \\ &:= T_1 + T_2 \end{aligned} \tag{9.25}$$

To prove the theorem, we will simply bound the supremum  $\sup_{f \in \mathcal{F}} |E(f)| \leq \sup_{f \in \mathcal{F}} |T_1 + T_2|$ . Pick a  $\delta$  covering  $\mathcal{F}_\delta$  of  $\mathcal{F}$  with size  $s \log \frac{C}{\delta}$  where  $\delta$  will be determined later in this proof. We first bound  $E(f)$  for all  $f \in \mathcal{F}_\delta$ . Given a fixed  $f$ , observe that  $\|\mathbf{z}\mathbf{m}(\mathbf{z}\mathbf{m}(y_i - f(\mathbf{x}_i))^2)\|_{\psi_1} \leq \mathcal{O}(\bar{K})^2 = \mathcal{O}(K + R)^2$  which follows from  $\|\mathbf{z}\mathbf{m}(y_i - f(\mathbf{x}_i))\|_{\psi_2} \leq \mathcal{O}(\bar{K}) = \mathcal{O}(K + R)$ . Applying subexponential concentration, since  $T_1$  is sum of i.i.d. subexponentials, we have

$$\mathbb{P}(|T_1| \geq \varepsilon) \leq \exp(-\mathcal{O}(n \min\{\varepsilon^2/\bar{K}^2, \varepsilon/\bar{K}\})). \tag{9.26}$$

Next, since  $\|\mathbf{z}\mathbf{m}(y_i - f(\mathbf{x}_i))\|_{\psi_2} \leq \mathcal{O}(\bar{K})$ , we can conclude that  $\|\frac{1}{n} \sum_{i=1}^n \mathbf{z}\mathbf{m}(y_i - f(\mathbf{x}_i))\|_{\psi_2} \leq \mathcal{O}(\bar{K})/\sqrt{n} \implies \|T_2\|_{\psi_1} \leq \mathcal{O}(\bar{K})^2/n$ . Using (9.26) for  $T_1$  and the subexponential tail bound for  $T_2$  holds when  $\varepsilon \leq \bar{K}$ , and assuming the number of samples  $n$  obeys  $n \geq \mathcal{O}(\max\{\varepsilon^{-1}, \varepsilon^{-2}\} \bar{K}^2 s \log \frac{C}{\delta})$ , then for all cover elements

$$|T_1| + |T_2| \leq 2\varepsilon.$$

holds with probability at least  $1 - \exp(-\mathcal{O}(n))$ . To conclude the proof, we need to move from the cover  $\mathcal{F}_\delta$  to  $\mathcal{F}$ . Pick  $f \in \mathcal{F}$  and its  $\delta$  neighbor  $f_\delta \in \mathcal{F}_\delta$ . Utilizing the deterministic relation  $|\mathbf{z}\mathbf{m}(X)| \leq |X| + |\mathbb{E}[X]|$  and using the fact that  $f_\delta$  is in a  $\delta$  neighborhood of  $f$ , we arrive at the following bounds

$$|\mathbf{z}\mathbf{m}(f(\mathbf{x}) - f_\delta(\mathbf{x}))| \leq \delta(\|\mathbf{x}\|_{\ell_2} + \mathbb{E}[\|\mathbf{x}\|_{\ell_2}]). \quad (9.27)$$

$$|\mathbf{z}\mathbf{m}(f(\mathbf{x}) + f_\delta(\mathbf{x}) - 2y)| \leq 2R(\|\mathbf{x}\|_{\ell_2} + \mathbb{E}[\|\mathbf{x}\|_{\ell_2}]) + 2K|\mathbf{z}\mathbf{m}(y)|. \quad (9.28)$$

Next observe that, with probability at least  $1 - 4n \exp(-p)$ , all  $\mathbf{x}_i, y_i$  obey  $\|\mathbf{x}_i\|_{\ell_2} \leq \mathcal{O}(\sqrt{p})$ ,  $|\mathbf{z}\mathbf{m}(y_i)| \leq \mathcal{O}(K\sqrt{p})$ . Combining this with (9.27), we conclude that for all  $1 \leq i \leq n$

$$|\mathbf{z}\mathbf{m}(f(\mathbf{x}_i) - y_i)^2 - \mathbf{z}\mathbf{m}(f_\delta(\mathbf{x}_i) - y_i)^2| \leq \mathcal{O}(\bar{K}\delta p). \quad (9.29)$$

Expanding the square differences in the same way, an identical argument shows the following two deviation bounds,

$$|\mathbb{E}[\mathbf{z}\mathbf{m}(f(\mathbf{x}_i) - y_i)^2 - \mathbf{z}\mathbf{m}(f_\delta(\mathbf{x}_i) - y_i)^2]| \leq \mathcal{O}(\bar{K}\delta p), \quad (9.30)$$

$$\frac{1}{n^2} \left| \left( \sum_{i=1}^n \mathbf{z}\mathbf{m}(y_i - f(\mathbf{x}_i)) \right)^2 - \left( \sum_{i=1}^n \mathbf{z}\mathbf{m}(y_i - f_\delta(\mathbf{x}_i)) \right)^2 \right| \leq \mathcal{O}(\bar{K}\delta p).$$

Combining these three inequalities ((9.29) and (9.30)) and substituting them in (9.25), we conclude that for all neighbors  $f_\delta, f$ ,

$$|E(f) - E(f_\delta)| \leq \mathcal{O}(\bar{K}\delta p).$$

Next we set  $\delta = c\varepsilon/(p\bar{K})$  for a sufficiently small constant  $c > 0$ , to find that with probability at least  $1 - \exp(-n)$ ,  $\sup_{f \in \mathcal{F}} |E(f)| \leq \varepsilon$  holds as long as the number of samples  $n$  obeys  $n \geq \mathcal{O}(\max\{\varepsilon^{-1}, \varepsilon^{-2}\} \bar{K}^2 s \log \frac{Cp\bar{K}}{c\varepsilon})$ . We define  $\mathcal{L}_{erm}(f) = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i) - r_{\text{avg}}(f))^2$  and  $\mathcal{L}_{pop}(f) = \mathbb{E}[\mathbf{z}\mathbf{m}(f(\mathbf{x}) - y)^2]$ . We also denote the CERM minimizer  $f_{erm} = \arg \min_{f \in \mathcal{F}} \mathcal{L}(f)$  and population minimizer  $f_{pop} = \min_{f \in \mathcal{F}} \mathcal{L}_{pop}(f)$ . Inequality (9.24) follows from the facts that we simultaneously have  $|E(f_{erm})| \leq \mathcal{O}(\varepsilon)$  and  $|E(f_{pop})| \leq \mathcal{O}(\varepsilon)$  which implies that

$$\mathcal{L}_{pop}(f_{erm}) \leq \mathcal{L}_{erm}(f_{erm}) + \mathcal{O}(\varepsilon) \leq \mathcal{L}_{erm}(f_{pop}) + \mathcal{O}(\varepsilon) \leq \mathcal{L}_{pop}(f_{pop}) + \mathcal{O}(\varepsilon),$$

concluding the proof. ■

## 9.6 Proof of Theorem 5.1

In this section we will show how Theorem 5.1 follows from Theorem 9.17. To this aim we need to show that  $\mathcal{F}_{\hat{\mathbf{k}}, B}$  has a small Lipschitz covering number. We construct the following cover  $\mathcal{F}'$  for the set  $\mathcal{F}_{\hat{\mathbf{k}}, B}$ . Let  $B' = B^{1/D}$ . Pick a  $\delta \leq B'$   $\ell_2$  cover  $\mathcal{C}$  of the interval  $[-B', B']$  which has size  $2B'/\delta$ . Let  $\mathcal{C}_i$  be identical copies of  $\mathcal{C}$ . We set

$$\mathcal{F}' = \{f_{\text{CNN}}(\beta_\ell \hat{\mathbf{k}}^{(\ell)}) \mid \beta_\ell \in \mathcal{C}_\ell, 1 \leq \ell \leq D\}.$$

In words, we construct CNNs by picking numbers from the cartesian product  $\mathcal{C}_1 \times \dots \times \mathcal{C}_D$  and scaling  $\{\hat{\mathbf{k}}^{(\ell)}\}_{\ell=1}^D$  with them. We now argue that  $\mathcal{F}'$  provides a cover of  $\mathcal{F}$ . Given  $f \in \mathcal{F}$  with scalings  $\beta_\ell$ , there exists  $f' \in \mathcal{F}'$  which uses scalings  $\beta'_\ell$  such that  $|\beta_\ell - \beta'_\ell| \leq \delta$ . Now, let  $f_\ell$  be the function with scalings  $\beta'_i$  until  $i = \ell$  and  $\beta_i$  for  $i > \ell$ . Note that  $f_0 = f$ ,  $f_D = f'$ . With this, we write

$$\|f - f'\|_{L_\infty} \leq \sum_{i=1}^D \|f_{i+1} - f_i\|_{L_\infty}.$$

Observe that  $f_{i-1}$  and  $f_i$  have equal layers except the  $i$ th layer. Let  $g_1$  be the function of the first  $i-1$  layers and  $g_2$  be the function of layers  $i+1$  to  $D$ . We have that  $f_{i+1}(\mathbf{x}) - f_i(\mathbf{x}) = g_2(\phi(\mathbf{K}_i(g_1(\mathbf{x})))) - g_2(\phi(\mathbf{K}'_i(g_1(\mathbf{x}))))$  where  $\mathbf{K}_i, \mathbf{K}'_i$  differ in the  $i$ th layer kernels of  $f$  and  $f'$  created from  $\beta_i \hat{\mathbf{k}}^{(i)}$  and  $\beta'_i \hat{\mathbf{k}}^{(i)}$  respectively. Also, observe that  $g_1$  is  $B^{i-1}$  Lipschitz and  $g_2(\phi(\cdot))$  is  $B^{D-i}$  Lipschitz functions. Hence,

$$|g_2(\phi(\mathbf{K}_i(g_1(\mathbf{x})))) - g_2(\phi(\mathbf{K}'_i(g_1(\mathbf{x}))))| \leq B^{D-i} |\mathbf{K}_i(g_1(\mathbf{x})) - \mathbf{K}'_i(g_1(\mathbf{x}))| \leq B^{D-i} \delta B^{i-1} \leq \delta B^{D-1}. \quad (9.31)$$

Summing over all  $i$ , this implies that  $\|f - f'\|_{L_\infty} \leq D\delta B^{D-1}$ . Recalling  $|\mathcal{F}'| \leq (2B'/\delta)^D$  and setting  $\delta = \varepsilon/(DB^{D-1})$ , the  $\varepsilon$  covering number of  $\mathcal{F}_{\hat{\mathbf{k}}, B}$  is  $N_\varepsilon \leq (2DB^{D-1}/\varepsilon)^D = (2DB/\varepsilon)^D$  which implies  $\log N_\varepsilon = D \log(\frac{2DB}{\varepsilon})$ . Now, since all kernels have Euclidean norm bounded by  $B'$ , we have  $\|f_{\text{CNN}}(\cdot)\|_{L_\infty} \leq B$  and  $\|f\|_{L_\infty} \leq B$  for all  $f \in \mathcal{F}$ . This also implies  $\|\mathbf{z}\mathbf{m}(f_{\text{CNN}}(\mathbf{x}))\|_{\psi_2} = \mathcal{O}(B)$ . Hence, we can apply Theorem 9.17 to conclude the proof of Theorem 5.1.

## Acknowledgements

M.S. is supported by the Air Force Office of Scientific Research Young Investigator Program (AFOSR-YIP) under award number FA9550-18-1-0078, a Google Faculty Research Award and a grant by the Northrop Grumman Cybersecurity Research Consortium.

## References

- [1] A. Alon Brutzkus and Amir Globerson. Globally optimal gradient descent for a convnet with Gaussian inputs. *arXiv preprint arXiv:1702.07966*, 2017.
- [2] Animashree Anandkumar, Rong Ge, Daniel Hsu, Sham M Kakade, and Matus Telgarsky. Tensor decompositions for learning latent variable models. *The Journal of Machine Learning Research*, 15(1):2773–2832, 2014.
- [3] Animashree Anandkumar, Rong Ge, and Majid Janzamin. Guaranteed non-orthogonal tensor decomposition via alternating rank-1 updates. *arXiv preprint arXiv:1402.5180*, 2014.
- [4] Sanjeev Arora, Aditya Bhaskara, Rong Ge, and Tengyu Ma. Provable bounds for learning some deep representations. In *International Conference on Machine Learning*, pages 584–592, 2014.
- [5] Rasmus Bro. Parafac. tutorial and applications. *Chemometrics and intelligent laboratory systems*, 38(2):149–171, 1997.
- [6] Emmanuel J. Candes, Xiaodong Li, and Mahdi Soltanolkotabi. Phase retrieval via wirtinger flow: Theory and algorithms. *arXiv preprint arXiv:1407.1065*, 2014.
- [7] Nadav Cohen, Or Sharir, and Amnon Shashua. On the expressive power of deep learning: A tensor analysis. In *Conference on Learning Theory*, pages 698–728, 2016.
- [8] Nadav Cohen and Amnon Shashua. Convolutional rectifier networks as generalized tensor decompositions. In *International Conference on Machine Learning*, pages 955–963, 2016.
- [9] Ronan Collobert and Jason Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167. ACM, 2008.
- [10] Sjoerd Dirksen. Tail bounds via generic chaining. *arXiv preprint arXiv:1309.3522*, 2013.
- [11] Simon S Du, Jason D Lee, and Yuandong Tian. When is a convolutional filter easy to learn? *arXiv preprint arXiv:1709.06129*, 2017.
- [12] Simon S Du, Jason D Lee, Yuandong Tian, Barnabas Poczos, and Aarti Singh. Gradient descent learns one-hidden-layer cnn: Don’t be afraid of spurious local minima. *arXiv preprint arXiv:1712.00779*, 2017.
- [13] Haoyu Fu, Yuejie Chi, and Yingbin Liang. Local geometry of one-hidden-layer neural networks for logistic regression. *arXiv preprint arXiv:1802.06463*, 2018.
- [14] Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points – online stochastic gradient for tensor decomposition. In *Conference on Learning Theory*, pages 797–842, 2015.

- [15] Rong Ge, Jason D Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. *arXiv preprint arXiv:1711.00501*, 2017.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [17] Kishore Jaganathan, Samet Oymak, and Babak Hassibi. Recovery of sparse 1-d signals from the magnitudes of their fourier transform. In *Information Theory Proceedings (ISIT), 2012 IEEE International Symposium On*, pages 1473–1477. IEEE, 2012.
- [18] Kishore Jaganathan et al. Predicting splicing from primary sequence with deep learning. *preprint*, 2018.
- [19] Majid Janzamin, Hanie Sedghi, and Anima Anandkumar. Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods. *arXiv preprint arXiv:1506.08473*, 2015.
- [20] Jean Kossaifi, Zachary C Lipton, Aran Khanna, Tommaso Furlanello, and Anima Anandkumar. Tensor regression networks. *arXiv preprint arXiv:1707.08308*, 2017.
- [21] Jean Kossaifi, Yannis Panagakis, and Maja Pantic. Tensorly: Tensor learning in python. *arXiv preprint arXiv:1610.09555*, 2016.
- [22] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [23] Yuanzhi Li and Yang Yuan. Convergence analysis of two-layer neural networks with relu activation. In *Advances in Neural Information Processing Systems*, pages 597–607, 2017.
- [24] Eran Malach and Shai Shalev-Shwartz. A provably correct algorithm for deep learning that actually works. *arXiv preprint arXiv:1803.09522*, 2018.
- [25] Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layers neural networks. *arXiv preprint arXiv:1804.06561*, 2018.
- [26] Marco Mondelli and Andrea Montanari. On the connection between learning two-layers neural networks and tensor decomposition. *arXiv preprint arXiv:1802.07301*, 2018.
- [27] Samet Oymak. Learning compact neural networks with regularization. *arXiv preprint arXiv:1802.01223*, 2018.
- [28] Samet Oymak and Mahdi Soltanolkotabi. Learning the input layer of a deep convolutional neural network via centered gradient descent. *preprint*, 2018.
- [29] Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. *arXiv preprint arXiv:1606.05336*, 2016.
- [30] Levent Sagun, Utku Evci, V Ugur Guney, Yann Dauphin, and Leon Bottou. Empirical analysis of the hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.
- [31] Nicholas D Sidiropoulos, Lieven De Lathauwer, Xiao Fu, Kejun Huang, Evangelos E Papalexakis, and Christos Faloutsos. Tensor decomposition for signal processing and machine learning. *IEEE Transactions on Signal Processing*, 65(13):3551–3582, 2017.
- [32] Mahdi Soltanolkotabi. Learning ReLUs via gradient descent. *arXiv preprint arXiv:1705.04591*, 2017.
- [33] Mahdi Soltanolkotabi. Structured signal recovery from quadratic measurements: Breaking sample complexity barriers via nonconvex optimization. *arXiv preprint arXiv:1702.06175*, 2017.
- [34] Mahdi Soltanolkotabi, Adel Javanmard, and Jason D Lee. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *arXiv preprint arXiv:1707.04926*, 2017.
- [35] Daniel Soudry and Yair Carmon. No bad local minima: Data independent training error guarantees for multilayer neural networks. *arXiv preprint arXiv:1605.08361*, 2016.
- [36] Michel Talagrand. *The generic chaining: upper and lower bounds of stochastic processes*. Springer Science & Business Media, 2006.
- [37] Michel Talagrand. Gaussian processes and the generic chaining. In *Upper and Lower Bounds for Stochastic Processes*, pages 13–73. Springer, 2014.
- [38] Ryota Tomioka and Taiji Suzuki. Spectral norm of random tensors. *arXiv preprint arXiv:1407.1870*, 2014.
- [39] Aaron Van den Oord, Sander Dieleman, and Benjamin Schrauwen. Deep content-based music recommendation. In *Advances in neural information processing systems*, pages 2643–2651, 2013.

- [40] Kai Zhong, Zhao Song, and Inderjit S Dhillon. Learning non-overlapping convolutional neural networks with multiple kernels. *arXiv preprint arXiv:1711.03440*, 2017.
- [41] Kai Zhong, Zhao Song, Prateek Jain, Peter L Bartlett, and Inderjit S Dhillon. Recovery guarantees for one-hidden-layer neural networks. *arXiv preprint arXiv:1706.03175*, 2017.