# How Many Samples are Needed to Estimate a Convolutional Neural Network?

Simon S. Du\* Carnegie Mellon University

Xiyu Zhai Massachusetts Institute of Technology

**Ruslan Salakhutdinov** Carnegie Mellon University Yining Wang\* Carnegie Mellon University

> Sivaraman Balakrishnan Carnegie Mellon University

Aarti Singh Carnegie Mellon University

# Abstract

A widespread folklore for explaining the success of Convolutional Neural Networks (CNNs) is that CNNs use a more compact representation than the Fullyconnected Neural Network (FNN) and thus require fewer training samples to accurately estimate their parameters. We initiate the study of rigorously characterizing the sample complexity of estimating CNNs. We show that for an *m*-dimensional convolutional filter with linear activation acting on a d-dimensional input, the sample complexity of achieving population prediction error of  $\epsilon$  is  $\tilde{O}(m/\epsilon^2)^2$ , whereas the sample-complexity for its FNN counterpart is lower bounded by  $\Omega(d/\epsilon^2)$  samples. Since, in typical settings  $m \ll d$ , this result demonstrates the advantage of using a CNN. We further consider the sample complexity of estimating a onehidden-layer CNN with linear activation where both the *m*-dimensional convolutional filter and the r-dimensional output weights are unknown. For this model, we show that the sample complexity is  $\widetilde{O}((m+r)/\epsilon^2)$  when the ratio between the stride size and the filter size is a constant. For both models, we also present lower bounds showing our sample complexities are tight up to logarithmic factors. Our main tools for deriving these results are a localized empirical process analysis and a new lemma characterizing the convolutional structure. We believe that these tools may inspire further developments in understanding CNNs.

# 1 Introduction

Convolutional Neural Networks (CNNs) have achieved remarkable impact in many machine learning applications, including computer vision (Krizhevsky et al., 2012), natural language processing (Yu et al., 2018) and reinforcement learning (Silver et al., 2016). The key building block of these improvements is the use of convolutional (weight sharing) layers to replace traditional fully connected layers, dating back to LeCun et al. (1995). A common folklore of explaining the success of CNNs is that they are a more compact representation than Fully-connected Neural Networks (FNNs) and thus require fewer samples to estimate. However, to our knowledge, there is no rigorous characterization of the sample complexity of learning a CNN.

The main difficulty lies in the convolution structure. Consider the simplest CNN, a single convolutional filter with linear activation followed by average pooling (see Figure 1a), which represents a

32nd Conference on Neural Information Processing Systems (NeurIPS 2018), Montréal, Canada.

<sup>\*</sup>Equal contribution.

<sup>&</sup>lt;sup>2</sup>We use the standard big-O notation in this paper and use  $\tilde{O}(\cdot)$  when we ignore poly-logarithmic factors.

function  $F_1 : \mathbb{R}^d \mapsto \mathbb{R}$  of the form:

$$F_1(x;w) = \sum_{\ell=0}^{r-1} w^{\top} \mathsf{P}_s^{\ell} x,$$
(1)

where  $w \in \mathbb{R}^m$  is the filter of size m and a stride size of  $s, r \approx d/s$  is the total number of times filter w is applied to an input vector  $x \in \mathbb{R}^d$ , and  $\mathsf{P}_s^\ell x := [x_{\ell s+1}, x_{\ell s+2}, \ldots, x_{\ell s+m}]$  is an m-dimensional segment of the feature vector x. Noting that  $F_1$  is a linear function of x, we can also represent  $F_1$  by a one-layer fully connected neural network (linear predictor):

$$F_1^{\text{FNN}}(x,\theta) = \theta^{\top} x \tag{2}$$

for some  $\theta \in \mathbb{R}^d$ . Suppose we have *n* samples  $\{x_i, y_i\}_{i=1}^n$  where *x* is the input and *y* is the label and use the least squares estimator:

$$\widehat{\theta} := \underset{\theta \in \mathbb{R}^d}{\operatorname{arg\,min}} \sum_{i=1}^n (y_i - \theta^\top x_i)^2.$$

By a classical results analyzing the prediction error for linear regression (see for instance (Wasserman, 2013)), under mild regularity conditions, we need  $n = d/\epsilon^2$  to have  $\sqrt{\mathbb{E}_{x \sim \mu} |\hat{\theta}^\top x - \theta_0^\top x|^2} \leq \epsilon$ , where  $\mu$  is the input distribution and  $\theta_0$  is the optimal linear predictor. The proof for FNN is fairly simple because we can write  $\hat{\theta} = (X^\top X)^{-1} X^\top Y$  (normal equation) where X and Y are the aggregated features and labels, respectively and then directly analyze this expression.

On the other hand, the network  $F_1$  can be viewed as a linear regression model with respect to w, by considering a "stacked" version of feature vectors  $\tilde{x}_i = \sum_{\ell=0}^{r-1} \mathsf{P}_s^\ell x \in \mathbb{R}^m$ . The classical analysis of ordinary least squares in linear regression does not directly yield the optimal sample complexity in this case, because the distributional properties of  $\tilde{x}_i$  as well as the spectral properties of the sample covariance  $\sum_i \tilde{x}_i \tilde{x}_i^\top$  are difficult to analyze due to the heavy correlation between coordinates of  $\tilde{x}$  corresponding to overlapping patches. We discuss further details of this aspect after our main positive result in Theorem 1.

In this paper, we take a step towards understanding the statistical behavior of the CNN model described above. We adopt tools from localized empirical process theory (van de Geer, 2000) and combine them with a structural property of convolutional filters (see Lemma 2) to give a complete characterization of the statistical behavior of this simple CNN.

We first consider the problem of learning a convolutional filter with average pooling as in Eq.(1) using the least squares estimator. We show in the standard statistical learning setting, under fairly natural conditions on the input distribution,  $\hat{w}$  satisfies

$$\sqrt{\mathbb{E}_{x\sim\mu}|F_1(x,\hat{w}) - F_1(x,w_0)|^2} = \widetilde{O}\left(\sqrt{m/n}\right),\,$$

where  $\mu$  is the input distribution and  $w_0$  is the underlying true convolutional filter. Notably, to achieve an  $\epsilon$  error, the CNN only needs  $\tilde{O}(m/\epsilon^2)$  samples whereas the FNN needs  $\Omega(d/\epsilon^2)$ . Since the filter size  $m \ll d$ , this result clearly justifies the folklore that the convolutional layer is a more compact representation. Furthermore, we complement this upper bound with a minimax lower bound which shows the error bound  $\tilde{O}(\sqrt{m/n})$  is tight up to logarithmic factors.

Next, we consider a one-hidden-layer CNN (see Figure 1b):

$$F_2(x; w, a) = \sum_{\ell=0}^{r-1} a_\ell w^\top \mathsf{P}_s^\ell x,$$
(3)

where both the shared convolutional filter  $w \in \mathbb{R}^m$  and output weights  $a \in \mathbb{R}^r$  are unknown. This architecture is previously considered in Du et al. (2017b). However the focus of that work is to understand the dynamics of gradient descent. Using similar tools as in analyzing a single convolutional filter, we show that the least squares estimator achieves the error bound  $\tilde{O}(\sqrt{(m+r)/n})$  if the ratio between the stride size and the filter size is a constant. Further, we present a minimax lower bound showing that the obtain rate is tight up to logarithmic factors.

To our knowledge, these theoretical results are the first sharp analyses of the statistical efficiency of the CNN. These results suggest that if the input follows a (linear) CNN model, then it can be learned more easily than treating it as a FNN since a CNN model reuses weights.





(a) Prediction function formalized in Eq. (1). It consists of a convolutional filter followed by averaged pooling. The convolutional filter is unknown.

(b) Prediction function formalized in Eq. (3) It consists of a convolutional filter followed by a linear prediction layer. Both layers are unknown.

Figure 1: CNN architectures that we consider in this paper.

### 1.1 Comparison with Existing Work

Our work is closely related to the analysis of the generalization ability of neural networks (Arora et al., 2018; Anthony & Bartlett, 2009; Bartlett et al., 2017b,a; Neyshabur et al., 2017; Konstantinos et al., 2017). These generalization bounds are often of the form:

$$L(\theta) - L_{\rm tr}(\theta) \leqslant D/\sqrt{n} \tag{4}$$

where  $\theta$  represents the parameters of a neural network,  $L(\cdot)$  and  $L_{tr}(\cdot)$  represent population and empirical error under some *additive* loss, and D is the model capacity and is finite only if the (spectral) norm of the weight matrix for each layer is bounded. Comparing with generalization bounds based on model capacity, our result has two advantages:

- 1. If  $L(\cdot)$  is taken to be the mean-squared<sup>3</sup> error  $\mathbb{E}|\cdot|^2$ , Eq. (4) implies an  $\widetilde{O}(1/\epsilon^4)$  sample complexity to achieve a standardized mean-square error of  $\sqrt{\mathbb{E}|\cdot|^2} \leq \epsilon$ , which is considerably larger than the  $\widetilde{O}(1/\epsilon^2)$  sample complexity we established in this paper.
- 2. Since the complexity of a model class in regression problems typically depends on the magnitude of model parameters (e.g.,  $||w||_2$ ), generalization error bounds like Eq. (4) are not scaleindependent and deteriorate if  $||w||_2$  is large. In contrast, our analysis has no dependency on the scale of w and also places no constraints on  $||w||_2$ .

On the other hand, we consider the special case where the neural network model is well-specified and the label is generated according to a neural network with unbiased additive noise (see Eq. (5)) whereas the generalization bounds discussed in this section are typically model agnostic.

#### 1.2 Other Related Work

Recently, researchers have been making progress in theoretically understanding various aspects of neural networks, including hardness of learning (Goel et al., 2016; Song et al., 2017; Brutzkus & Globerson, 2017), landscape of the loss function (Kawaguchi, 2016; Choromanska et al., 2015; Hardt & Ma, 2016; Haeffele & Vidal, 2015; Freeman & Bruna, 2016; Safran & Shamir, 2016; Zhou & Feng, 2017; Nguyen & Hein, 2017b,a; Ge et al., 2017b; Zhou & Feng, 2017; Safran & Shamir, 2017; Du & Lee, 2018), dynamics of gradient descent (Tian, 2017; Zhong et al., 2017b; Li & Yuan, 2017), provable learning algorithms (Goel & Klivans, 2017a,b; Zhang et al., 2015), etc.

Focusing on the convolutional neural network, most existing work has analyzed the convergence rate of gradient descent or its variants (Du et al., 2017a,b; Goel et al., 2018; Brutzkus & Globerson, 2017; Zhong et al., 2017a). Our paper differs from them in that we do not consider the computational complexity but only the sample complexity and information theoretical limits of learning a CNN. It is an open question when taking computational budget into account, what is the optimal estimator for CNN.

<sup>&</sup>lt;sup>3</sup>Because the standardized mean-square error  $\sqrt{\mathbb{E}|\cdot|^2}$  is not a sum of independent random variables, it is difficult, if not impossible, to apply generalization error bounds directly for  $\sqrt{\mathbb{E}|\cdot|^2}$ .

Convolutional structure has also been studied in the dictionary learning (Singh et al., 2018; Huang & Anandkumar, 2015) and blind de-convolution (Zhang et al., 2017) literature. These papers studied the unsupervised setting where their goal is to recover structured signals from observations generated according to convolution operations whereas our paper focuses on the supervised learning setting with predictor having the convolution structure.

#### 1.3 Organization

This paper is organized as follows. In Section 2, we formally setup the problem and assumptions. In Section 3 we present our main theoretical results for learning a convolutional filter (see Eq. (1)). In Section 4 we present our main theoretical results for learning a one-hidden-layer CNN (see Eq. (3)). In Section 5, we use numerical experiments to verify our theoretical findings. We conclude and list future directions in Section 6. Most technical proofs are deferred to the appendix.

## **2** Problem specification and assumptions

Let  $\{x_i, y_i\}_{i=1}^n$  be a sample of *n* training data points, where  $x_i \in \mathbb{R}^d$  denotes the *d*-dimensional feature vector of the *i*th data point and  $y_i \in \mathbb{R}$  is the corresponding real-valued response. We consider a generic model of

$$y_i = F(x_i; \boldsymbol{w}_0) + \varepsilon_i, \quad \text{where } \mathbb{E}[\varepsilon_i | x_i] = 0.$$
 (5)

In the model of Eq. (5), F represents a certain network parameterized by a fixed but unknown parameter  $w_0$  that takes a *d*-dimensional vector  $x_i$  as input and outputs a single real-valued prediction  $F(x_i; w_0)$ .  $\{\varepsilon_i\}_{i=1}^n$  represents stochastic noise inherent in the data, and is assumed to have mean zero. The feature vectors of training data  $\{x_i\}_{i=1}^n$  are sampled i.i.d. from an unknown distribution  $\mu$  supported on  $\mathbb{R}^d$ .

Throughout this paper we make the following assumptions:

- (A1) Sub-gaussian noise: there exists constant  $\sigma^2 < \infty$  such that for any  $t \in \mathbb{R}$ ,  $\mathbb{E}e^{t\varepsilon_i} \leq e^{\sigma^2 t^2/2}$ ;
- (A2) Sub-gaussian design: there exists constant  $\nu^2 < \infty$  such that for any  $a \in \mathbb{R}^d$ ,  $\mathbb{E}_{\mu} x = 0$  and  $\mathbb{E}_{\mu} \exp\{a^{\top} x\} \leq \exp\{\nu^2 \|a\|_2^2/2\};$
- (A3) Non-degeneracy: there exists constant  $\kappa > 0$  such that  $\lambda_{\min}(\mathbb{E}_{\mu}xx^{\top}) \ge \kappa$ .

We remark that the assumptions (A1) through (A3) are quite mild. In particular, we only impose sub-Gaussianity conditions on the distributions of  $x_i$  and  $\varepsilon_i$ , and do not assume they are generated/sampled from any *exact* distributions. The last non-degeneracy condition (A3) assumes that there is a non-negligible probability mass along any direction of the input distributions. It is very likely to be satisfied after simple pre-processing steps of input data, such as mean removal and whitening of the sample covariance.

We are interested in learning a parameter  $\hat{w}_n$  using a training sample  $\{(x_i, y_i)\}_{i=1}^n$  of size n so as to minimize the standardized *population* mean-square prediction error

$$\operatorname{err}_{\mu}(\widehat{\boldsymbol{w}}_{n}, \boldsymbol{w}_{0}; F) = \sqrt{\mathbb{E}_{x \sim \mu} \left| F(x; \widehat{\boldsymbol{w}}_{n}) - F(x; \boldsymbol{w}_{0}) \right|^{2}}.$$
(6)

# **3** Convolutional filters with average pooling

We first consider a convolutional network with one convolutional layer, one convolutional filter, an average pooling layer and linear activations. More specifically, for a single convolutional filter  $w \in \mathbb{R}^m$  of size m and a stride of size s, the network can be written as

$$F_1(x;w) = \sum_{\ell=0}^{r-1} w^{\top} \mathsf{P}_s^{\ell} x,$$
(7)

where  $r \approx d/s$  is the total number of times filter w is applied to an input vector x, and  $\mathbb{P}_s^{\ell} x := [x_{\ell s+1}, x_{\ell s+2}, \dots, x_{\ell s+m}]$  is an *m*-dimensional segment of the *d*-dimensional feature vector  $x_i$ . For simplicity, we assume that m is divisible s and let  $J = m/s \in \mathbb{N}$  denote the number of strides within a single filter of size m.

#### 3.1 The upper bound

Given training sample  $\{(x_i, y_i)\}_{i=1}^n$ , we consider the following least-squares estimator:

$$\widehat{w}_n \in \arg\min_{w \in \mathbb{R}^m} \frac{1}{n} \sum_{i=1}^n \left( y_i - F_1(x_i; w) \right)^2.$$
(8)

Note the subscript n which emphasizes that  $\hat{w}_n$  is trained using a sample of n data points. In addition, because the objective is a quadratic function in w, Eq. (8) is actually a convex optimization problem and a global optimal solution  $\hat{w}_n$  can be obtained efficiently. More specifically,  $\hat{w}_n$  admits the closed-form solution of  $\hat{w}_n = (\sum_{i=1}^n \tilde{x}_i \tilde{x}_i^\top)^{-1} \sum_{i=1}^n y_i \tilde{x}_i$ , where  $\tilde{x}_i = \sum_{\ell=0}^{r-1} \mathsf{P}_s^\ell x_i$  is the stacked version of input feature vector  $x_i$ .

The following theorem upper bounds the expected population mean-square prediction error  $\operatorname{err}_{\mu}(\hat{w}_n, w_0; F_1)$  of the least-square estimate  $\hat{w}_n$  in Eq. (8).

**Theorem 1.** Fix an arbitrary  $\delta \in (0, 1/2)$ . Suppose (A1) through (A3) hold and  $\nu \sqrt{\log(n/\delta)} \ge \kappa$ ,  $n \ge \kappa^{-2}\nu^2 m \log(\nu d \log \delta^{-1}) \log(n\delta^{-1})$ . Then there exists a universal constant C > 0 such that with probability  $1 - \delta$  over the random draws of  $x_1, \ldots, x_n \sim \mu$ ,

$$\mathbb{E}\mathrm{err}_{\mu}(\hat{w}_n, w_0; F_1) \leqslant C\sqrt{\frac{\sigma^2 m \log(\kappa^{-1}\nu d \log(\delta^{-1}))}{n}} \quad \text{conditioned on } x_1, \dots, x_n.$$
(9)

*Here the expectation is taken with respect to the randomness in*  $\{\varepsilon_i\}_{i=1}^n$ .

Theorem 1 shows that, with  $n = \tilde{\Omega}(m)$  samples, the expected population mean-square error  $\operatorname{err}_{\mu}(\hat{w}_n, w_0; F_1)$  scales as  $\tilde{O}(\sqrt{\sigma^2 m/n})$ . This matches the  $1/\sqrt{n}$  statistical error for classical parametric statistics problems, and also confirms the "parameter count" intution that the estimation error scales approximately with the number of parameters in a network  $(m \text{ in network } F_1)$ .

We next briefly explain the strategies we employ to prove Theorem 1. While it's tempting to directly use the closed-form expression  $\hat{w}_n = (\sum_{i=1}^n \tilde{x}_i \tilde{x}_i^{\top})^{-1} \sum_{i=1}^n y_i \tilde{x}_i$  to analyze  $\hat{w}_n$ , such an approach has two limitations. First, because we consider the *population* mean-square error  $\operatorname{err}_{\mu}(\hat{w}_n, w_0; F_1)$ , such an approach would inevtiably require the analysis of spectral properties (e.g., the least eigenvalue) of  $\sum_{i=1}^n \tilde{x}_i \tilde{x}_i^{\top}$ , which is very challenging as heavy correlation occurs in  $\tilde{x}_i$  when filters are overlapping (i.e., s < m and J > 1). It is likely that strong assumptions such as *exact isotropic* Gaussianity of the feature vectors are needed to analyze the distributional properties  $\tilde{x}_i$  (Qu et al., 2017). Also, such an approach relies on closed-forms of  $\hat{w}_n$  and is difficult to extend to other potential activations such as the ReLU activation. when no closed-form expressions of  $\hat{w}_n$  exist.

To overcome the above difficulties, we adopt a *localized empirical process* approach introduced in (van de Geer, 2000) to upper bound the expected population mean-square prediction error. At the core analysis is an upper bound on the covering number of a *localized parameter set*, with an interesting argument that partitions a *d*-dimensional equivalent regressor for compactification purposes (see Lemmas 2 and 4 in the appendix for details). Our proof does not rely on the exact/closed-form expression of  $\hat{w}_n$ , and has the potential to be extended to other activation functions, as we discuss in Section 6. The complete proof of Theorem 1 is placed in the appendix.

#### 3.2 The lower bound

We prove the following information-theoretic lower bound on  $\mathbb{E}\operatorname{err}_{\mu}(\widehat{w}_n, w_0)$  of any estimator  $\widehat{w}_n$  calculated on a training sample of size n.

**Theorem 2.** Suppose  $x_1, \ldots, x_n \sim \mathcal{N}(0, I)$  and  $\varepsilon_1, \ldots, \varepsilon_n \sim \mathcal{N}(0, \sigma^2)$ . Suppose also that m - s is an even number. Then there exists a universal constant C' > 0 such that

$$\inf_{\widehat{w}_n} \sup_{w_0 \in \mathbb{R}^m} \mathbb{E}\mathrm{err}_{\mu}(\boldsymbol{w}_n, w_0; F_1) \ge C' \sqrt{\frac{\sigma^2 m}{n}}.$$
(10)

**Remark 1.** Theorem 2 is valid for any pair of (filter size, stride) combinations (m, s), provided that m is divisible by s and m - s is an even number. The latter requirement is a technical condition in our proof and is not critical, because one can double the size of m and s, and the lower bound in Theorem 2 remains asymptotically on the same order.

Theorem 2 shows that any estimator  $\hat{w}_n$  computed on a training set of size n must have a worst-case error of at least  $\sqrt{\sigma^2 m/n}$ . This suggests that our upper error bound in Theorem 1 is tight up to logarithmic factors.

Our proof of Theorem 2 draws on tools from standard information-theoretical lower bounds such as the Fano's inequality (Yu, 1997; Tsybakov, 2009). The high-level idea is to construct a *finite* candidate set of parameters  $\mathcal{W} \subseteq \mathbb{R}^m$  and upper bound the Kullback-Leibler (KL) divergence of induced observable distributions and the population prediction mean-square error between parameters in the candidate set  $\mathcal{W}$ . The complete proof of Theorem 2 is placed in the appendix.

## 4 Convolutional filters with prediction layers

We consider a slightly more complicated convolutional network with two layers: the first layer is a single convolutional filter of size m, applied r times to a d-dimensional input vector with stride s; the second layer is a linear regression prediction layer that produces a single real-valued output.

For such a two-layer network the parameter w can be specified as w = (w, a), where  $w \in \mathbb{R}^m$  is the weights in the first-layer convolutional filter and  $a \in \mathbb{R}^r$  is the weight in the second linear prediction layer. The network  $F_2(x; w) = F_2(x; w, a)$  can then be written as

$$F_2(x; w, a) = \sum_{\ell=0}^{r-1} a_\ell w^\top \mathsf{P}_s^\ell x.$$
 (11)

Note that in Eq. (11) the vector  $a \in \mathbb{R}^r$  is labeled as  $a = (a_0, a_1, \dots, a_{r-1})$  for convenience that matches with the labels of the operator  $\mathsf{P}_s^\ell$  for  $\ell = 0, \dots, r-1$ .

Compared to network  $F_1$  with average pooling, the new network  $F_2$  can be viewed as a *weighted* pooling of convolutional filters, with weights  $a \in \mathbb{R}^r$  unknown and to be learnt. A graph illustration of the network  $F_2$  is given in Figure 1b.

#### 4.1 The upper bound

We again consider the least-squares estimator

$$\widehat{\boldsymbol{w}}_n = (\widehat{w}_n, \widehat{a}_n) \in \arg\min_{\boldsymbol{w} \in \mathbb{R}^m, \boldsymbol{a} \in \mathbb{R}^r} \frac{1}{n} \sum_{i=1}^n \left( y_i - F_2(x_i; \boldsymbol{w}, \boldsymbol{a}) \right)^2.$$
(12)

Again, we use subscript n to emphasize that both  $\hat{w}_n$  and  $\hat{a}_n$  are computed on a training set  $\{x_i, y_i\}_{i=1}^n$  of size n.

Unlike the least squares problem in Eq. (8) for the  $F_1$  network, the optimization problem in Eq. (12) has two optimization variables w, a and is therefore no longer convex. This means that popular optimization algorithms like gradient descent do not necessarily converge to a global minima in Eq. (12). Nevertheless, in this paper we choose to focus on the *statistical* properties of  $(\hat{w}_n, \hat{a}_n)$  and assume global minimality of Eq. (12) is achieved. On the other hand, because Eq. (12) resembles the matrix sensing problem, it is possible that all local minima are global minima and saddle points can be efficiently escaped (Ge et al., 2017a), which we leave as future work.

The following theorem upper bounds the population mean-square prediction error of any global minimizer  $\hat{w}_n = (\hat{w}_n, \hat{a}_n)$  of Eq. (12).

**Theorem 3.** Fix arbitrary  $\delta \in (0, 1/2)$  and define J := m/s, where m is the filter size and s is the stride. Suppose (A1) through (A3) hold and  $\nu \sqrt{\log(n/\delta)} \ge \kappa$ ,  $n \ge \kappa^{-2}\nu^2(rJ + m)\log(\nu d\log \delta^{-1})\log(n\delta^{-1})$ . Then there exists a universal constant C > 0 such that with probability  $1 - \delta$  over the random draws of  $x_1, \ldots, x_n \sim \mu$ ,

$$\mathbb{E}\mathrm{err}_{\mu}(\hat{\boldsymbol{w}}_n, \boldsymbol{w}_0; F_2) \leqslant C\sqrt{\frac{\sigma^2(rJ+m)\log(\kappa^{-1}\nu d\log(\delta^{-1}))}{n}} \quad conditioned \text{ on } x_1, \dots, x_n.$$
(13)

*Here the expectation is taken with respect to the randomness in*  $\{\varepsilon_i\}_{i=1}^n$ .

Theorem 3 is proved by a similar localized empirical process arguments as in the proof of Theorem 1. Due to space costraints we defer the complete proof of Theorem 3 to the appendix.



Figure 2: Experiments on the problem of learning a convolutional filter with average pooling described in Section 3 with stride size s = 1.

Theorem 3 shows that  $\operatorname{err}_{\mu}(\widehat{w}_n, w_0; F_2)$  can be upper bounded by  $\widetilde{O}(\sqrt{\sigma^2(rJ+m)/n})$ , provided that at least  $n = \widetilde{\Omega}(rJ+m)$  samples are available. Compared to the intuitive "parameter count" of r+m (r parameters for a and m parameters for w), our upper bound has an additional multiplicative J = m/s term, which is the number of strides within each m-dimensional filter. Therefore, our upper bound only matches parameter counts when J is very small (e.g., non-overlapping filters or fast-moving filters where the stride s is at least a constant fraction of filter size m), and becomes large when the stride s is very small, leading to many convolutions being computed.

We conjecture that such an increase in error/sample complexity is due to an inefficiency in one of our key technical lemmas. More specifically, in Lemma 7 in which we derive upper bounds on covering number of localized parameter sets, we use the boundedness and low-dimensionality of each segment of differences of equivalent parameters for compactification purposes; such an argument is not ideal, as it overlooks the correlation between different segments, connected by an r-dimensional parameter a. A sharper covering number argument would potentially improve the error analysis and achieve sample complexity scaling with r + m.

#### 4.2 The lower bound

We prove the following information-theoretical lower bound on  $\mathbb{E}\operatorname{err}_{\mu}(\hat{\boldsymbol{w}}_n, \boldsymbol{w}_0)$  of any estimator  $\hat{\boldsymbol{w}}_n = (\hat{w}_n, \hat{a}_n)$  calculated on a training sample of size n.

**Theorem 4.** Suppose  $x_1, \ldots, x_n \sim \mathcal{N}(0, I)$  and  $\varepsilon_1, \ldots, \varepsilon_n \sim \mathcal{N}(0, \sigma^2)$ . Then there exists a universal constant C' > 0 such that

$$\inf_{\hat{\boldsymbol{w}}_n} \sup_{\boldsymbol{w}_0} \mathbb{E}\mathrm{err}_{\mu}(\hat{\boldsymbol{w}}_n, \boldsymbol{w}_0; F_2) \ge C' \sqrt{\frac{\sigma^2(r+m)}{n}}.$$
(14)

Theorem 4 shows that the error of any estimator  $\hat{w}_n$  computed on a training sample of size n must scale as  $\sqrt{\sigma^2(r+m)/n}$ , matching the parameter counts of r+m for  $F_2$ . It is proved by reducing the regression problem under  $F_2$  to two separate ordinary linear regression problems and invoking classical lower bounds for linear regression models (Wasserman, 2013; Van der Vaart, 1998). A complete proof of Theorem 4 is given in the appendix.

# **5** Experiments

In this section we use simulations to verify our theoretical findings. For all experiments, we let the ambient dimension d be 64 and the input distribution be Gaussian with mean 0 and identity covariance. We use the population mean-square prediction error defined in Eq. (6) as the evaluation metric. In all plots, CNN represents using convolutional parameterization corresponding to Eq. (1) or Eq. (3) and FNN represents using fully connected parametrization corresponding to Eq. (2).

In Figure 2 and Figure 3, we consider the problem of learning a convolutional filter with average pooling which we analyzed in Section 3. We vary the number of samples, the dimension of filters and the stride size. Here we compare parameterizing the prediction function as a *d*-dimensional linear predictor and as a convolutional filter followed by average pooling. Experiments show CNN



Figure 3: Experiments on the problem of learning a convolutional filter with average pooling described in Section 3 with stride size s = m, i.e., non-overlapping.



Figure 4: Experiment on the problem of one-hidden-layer convolutional neural network with a shared filter and a prediction layer described in Section 4. The filter size m is chosen to be 8.

parameterization is consistently better than the FNN parameterization. Further, as number of training samples increases, the prediction error goes down and as the dimension of filter increases, the error goes up. These facts qualitatively justify our derived error bound  $\tilde{O}\left(\frac{m}{n}\right)$ . Lastly, in Figure 2 we choose stride s = 1 and in Figure 3 we choose stride size equals to the filter size s = m, i.e., non-overlapping. Our experiment shows the stride does *not* affect the prediction error in this setting which coincides our theoretical bound in which there is no stride size factor.

In Figure 4, we consider the one-hidden-layer CNN model analyzed in Section 4. Here we fix the filter size m = 8 and vary the number of training samples and the stride size. When stride s = 1, convolutional parameterization has the same order parameters as the linear predictor parameterization  $(r = 57 \text{ so } r + m = 65 \approx d = 64)$  and Figure 4a shows they have similar performances. In Figure 4b and Figure 4c we choose the stride to be m/2 = 4 and m = 8 (non-overlapping), respectively. Note these settings have less parameters (r + m = 23 for s = 4 and r + m = 16 for s = 8) than the case when s = 1 and so CNN gives better performance than FNN.

# 6 Conclusion and Future Directions

In this paper we give rigorous characterizations of the statistical efficiency of CNN with simple architectures. Now we discuss how to extend our work to more complex models and main difficulties.

**Non-linear Activation:** Our paper only considered CNN with linear activation. A natural question is what is the sample complexity of learning a CNN with non-linear activation like Recitifed Linear Units (ReLU). We find that even without convolution structure, this is a difficult problem. For linear activation function, we can show the empirical loss is a good approximation to the population loss and we used this property to derive our upper bound. However, for ReLU activation, we can find a counter example for any finite n, which breaks our Lemma 3. We believe if there is a better understanding of non-smooth activation which can replace our Lemma 3, we can extend our analysis framework to derive sharp sample complexity bounds for CNN with non-linear activation function.

**Multiple Filters:** For both models we considered in this paper, there is only one shared filter. In commonly used CNN architectures, there are multiple filters in each layer and multiple layers. Note

that if one considers a model of k filters with linear activation with k > 1, one can always replace this model by a single convolutional filter that equals to the summation of these k filters. Thus, we can formally study the statistical behavior of wide and deep architectures only after we have understood the non-linear activation function. Nevertheless, we believe our empirical process based analysis is still applicable.

# Acknowledgment

This research was partly funded by AFRL grant FA8750-17-2-0212 and DARPA D17AP00001.

## References

- Anthony, M., & Bartlett, P. L. (2009). Neural network learning: Theoretical foundations. cambridge university press.
- Arora, S., Ge, R., Neyshabur, B., & Zhang, Y. (2018). Stronger generalization bounds for deep nets via a compression approach. arXiv preprint arXiv:1802.05296.
- Bartlett, P. L., Foster, D. J., & Telgarsky, M. J. (2017a). Spectrally-normalized margin bounds for neural networks. In *Advances in Neural Information Processing Systems*, (pp. 6241–6250).
- Bartlett, P. L., Harvey, N., Liaw, C., & Mehrabian, A. (2017b). Nearly-tight vcdimension and pseudodimension bounds for piecewise linear neural networks. arxiv preprint. *arXiv*, 1703.
- Bickel, P. J., Ritov, Y., & Tsybakov, A. B. (2009). Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics*, 37(4), 1705–1732.
- Brutzkus, A., & Globerson, A. (2017). Globally optimal gradient descent for a Convnet with Gaussian inputs. arXiv preprint arXiv:1702.07966.
- Choromanska, A., Henaff, M., Mathieu, M., Arous, G. B., & LeCun, Y. (2015). The loss surfaces of multilayer networks. In *Artificial Intelligence and Statistics*, (pp. 192–204).
- Du, S. S., & Lee, J. D. (2018). On the power of over-parametrization in neural networks with quadratic activation. *arXiv preprint arXiv:1803.01206*.
- Du, S. S., Lee, J. D., & Tian, Y. (2017a). When is a convolutional filter easy to learn? *arXiv preprint arXiv:1709.06129*.
- Du, S. S., Lee, J. D., Tian, Y., Poczos, B., & Singh, A. (2017b). Gradient descent learns one-hiddenlayer cnn: Don't be afraid of spurious local minima. arXiv preprint arXiv:1712.00779.
- Dudley, R. M. (1967). The sizes of compact subsets of hilbert space and continuity of gaussian processes. *Journal of Functional Analysis*, 1(3), 290–330.
- Freeman, C. D., & Bruna, J. (2016). Topology and geometry of half-rectified network optimization. arXiv preprint arXiv:1611.01540.
- Ge, R., Jin, C., & Zheng, Y. (2017a). No spurious local minima in nonconvex low rank problems: A unified geometric analysis. In *Proceedings of the 34th International Conference on Machine Learning*, (pp. 1233–1242).
- Ge, R., Lee, J. D., & Ma, T. (2017b). Learning one-hidden-layer neural networks with landscape design. arXiv preprint arXiv:1711.00501.
- Goel, S., Kanade, V., Klivans, A., & Thaler, J. (2016). Reliably learning the ReLU in polynomial time. arXiv preprint arXiv:1611.10258.
- Goel, S., & Klivans, A. (2017a). Eigenvalue decay implies polynomial-time learnability for neural networks. *arXiv preprint arXiv:1708.03708*.
- Goel, S., & Klivans, A. (2017b). Learning depth-three neural networks in polynomial time. *arXiv* preprint arXiv:1709.06010.

- Goel, S., Klivans, A., & Meka, R. (2018). Learning one convolutional layer with overlapping patches. *arXiv preprint arXiv:1802.02547*.
- Graham, R., & Sloane, N. (1980). Lower bounds for constant weight codes. *IEEE Transactions on Information Theory*, 26(1), 37–43.
- Haeffele, B. D., & Vidal, R. (2015). Global optimality in tensor factorization, deep learning, and beyond. arXiv preprint arXiv:1506.07540.
- Hardt, M., & Ma, T. (2016). Identity matters in deep learning. arXiv preprint arXiv:1611.04231.
- Hoeffding, W. (1963). Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301), 13–30.
- Huang, F., & Anandkumar, A. (2015). Convolutional dictionary learning through tensor factorization. In *Feature Extraction: Modern Questions and Challenges*, (pp. 116–129).
- Kawaguchi, K. (2016). Deep learning without poor local minima. In Advances in Neural Information Processing Systems, (pp. 586–594).
- Konstantinos, P., Davies, M., & Vandergheynst, P. (2017). Pac-bayesian margin bounds for convolutional neural networks-technical report. arXiv preprint arXiv:1801.00171.
- Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In Advances in neural information processing systems, (pp. 1097–1105).
- LeCun, Y., Bengio, Y., et al. (1995). Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, *3361*(10), 1995.
- Li, Y., & Yuan, Y. (2017). Convergence analysis of two-layer neural networks with ReLU activation. arXiv preprint arXiv:1705.09886.
- Neyshabur, B., Bhojanapalli, S., McAllester, D., & Srebro, N. (2017). A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. *arXiv preprint arXiv:1707.09564*.
- Nguyen, Q., & Hein, M. (2017a). The loss surface and expressivity of deep convolutional neural networks. *arXiv preprint arXiv:1710.10928*.
- Nguyen, Q., & Hein, M. (2017b). The loss surface of deep and wide neural networks. *arXiv preprint arXiv:1704.08045*.
- Qu, Q., Zhang, Y., Eldar, Y. C., & Wright, J. (2017). Convolutional phase retrieval via gradient descent. arXiv preprint arXiv:1712.00716.
- Safran, I., & Shamir, O. (2016). On the quality of the initial basin in overspecified neural networks. In *International Conference on Machine Learning*, (pp. 774–782).
- Safran, I., & Shamir, O. (2017). Spurious local minima are common in two-layer relu neural networks. arXiv preprint arXiv:1712.08968.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. (2016). Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587), 484–489.
- Singh, S., Póczos, B., & Ma, J. (2018). Minimax reconstruction risk of convolutional sparse dictionary learning. In *International Conference on Artificial Intelligence and Statistics*, (pp. 1327– 1336).
- Song, L., Vempala, S., Wilmes, J., & Xie, B. (2017). On the complexity of learning neural networks. In *Advances in Neural Information Processing Systems*, (pp. 5520–5528).
- Tian, Y. (2017). An analytical formula of population gradient for two-layered ReLU network and its applications in convergence and critical point analysis. *arXiv preprint arXiv:1703.00560*.
- Tsybakov, A. B. (2009). *Introduction to nonparametric estimation*. Springer Series in Statistics. Springer, New York.

van de Geer, S. A. (2000). Empirical Processes in M-estimation, vol. 6. Cambridge university press.

- Van der Vaart, A. W. (1998). Asymptotic statistics, vol. 3. Cambridge university press.
- Vershynin, R. (2012). How close is the sample covariance matrix to the actual covariance matrix? *Journal of Theoretical Probability*, 25(3), 655–686.
- Wang, Y., & Singh, A. (2016). Noise-adaptive margin-based active learning and lower bounds under tsybakov noise condition. In *AAAI*.
- Wasserman, L. (2013). All of statistics: a concise course in statistical inference. Springer Science & Business Media.
- Yu, A. W., Dohan, D., Luong, M.-T., Zhao, R., Chen, K., Norouzi, M., & Le, Q. V. (2018). Qanet: Combining local convolution with global self-attention for reading comprehension. arXiv preprint arXiv:1804.09541.
- Yu, B. (1997). Assouad, fano, and le cam. In Festschrift for Lucien Le Cam, (pp. 423-435). Springer.
- Zhang, Y., Lau, Y., Kuo, H.-w., Cheung, S., Pasupathy, A., & Wright, J. (2017). On the global geometry of sphere-constrained sparse blind deconvolution. In *Proceedings of the IEEE Conference* on Computer Vision and Pattern Recognition, (pp. 4894–4902).
- Zhang, Y., Lee, J. D., Wainwright, M. J., & Jordan, M. I. (2015). Learning halfspaces and neural networks with random initialization. *arXiv preprint arXiv:1511.07948*.
- Zhong, K., Song, Z., & Dhillon, I. S. (2017a). Learning non-overlapping convolutional neural networks with multiple kernels. *arXiv preprint arXiv:1711.03440*.
- Zhong, K., Song, Z., Jain, P., Bartlett, P. L., & Dhillon, I. S. (2017b). Recovery guarantees for one-hidden-layer neural networks. arXiv preprint arXiv:1706.03175.
- Zhou, P., & Feng, J. (2017). The landscape of deep learning algorithms. *arXiv preprint* arXiv:1705.07038.

# **Appendix: Proofs**

We first define some notations that will be used throughout the proofs of our results. Note that because the activation function we consider in this paper is the identity mapping, both networks  $F_1$  and  $F_2$  can be written as a "structured linear regression" model

$$F_{1/2}(x; \boldsymbol{w}) = \langle x, \theta(\boldsymbol{w}) \rangle$$
 where  $\theta(\boldsymbol{w}) \in \Theta_{1/2} \subseteq \mathbb{R}^d$ . (15)

Here  $\Theta_{1/2}$  is a subset of  $\mathbb{R}^d$  subject to additional structural constraints corresponding to  $F_1$  or  $F_2$ . We can then define/rewrite "population" and "empirical" mean-square prediction errors as

$$\operatorname{err}_{\mu}^{2}(\boldsymbol{w},\boldsymbol{w}';F) = \mathbb{E}_{\mu} |\langle x,\theta(\boldsymbol{w})-\theta(\boldsymbol{w}')\rangle|^{2} =: \|\theta(\boldsymbol{w})-\theta(\boldsymbol{w}')\|_{\mu}^{2};$$
$$\operatorname{err}_{X}^{2}(\boldsymbol{w},\boldsymbol{w}';F) = \frac{1}{n} \sum_{i=1}^{n} |\langle x_{i},\theta(\boldsymbol{w})-\theta(\boldsymbol{w}')\rangle|^{2} =: \|\theta(\boldsymbol{w})-\theta(\boldsymbol{w}')\|_{X}^{2};$$

For any set  $\Theta \subseteq \mathbb{R}^d$ , error parameter  $\epsilon > 0$  and a distance metric  $d(\cdot, \cdot) : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}^+$ , define  $N(\epsilon; \Theta, d)$  as the *covering number* of  $\Theta$  in  $d(\cdot, \cdot)$ , which is the size of the smallest finite cover set  $H \subseteq \mathbb{R}^d$  such that  $\sup_{\phi \in \Theta} \min_{\phi' \in H} d(\phi, \phi') \leq \epsilon$ .

# A Proof of Theorem 1 (upper bound, average pooling)

For the convolutional model with an average pooling layer  $F_1(x; w) = \sum_{\ell=0}^{r-1} w^\top \mathsf{P}_s^\ell x$ , a linear model  $\theta(w) \in \mathbb{R}^d$  as in Eq. (15) can be produced as

$$\theta(w) = \sum_{\ell=0}^{r-1} \mathsf{S}_s^{\ell} w \quad \text{where } \; \mathsf{S}_s^{\ell} w = [\underbrace{0, \dots, 0}_{\ell s \; \text{zeros}}, w_1, \dots, w_m, 0, \dots, 0] \in \mathbb{R}^d.$$

Our first lemma is the following "basic inequality", which upper bounds  $\widehat{\operatorname{err}}_X^2(\widehat{w}_n, w_0; F) = \|\theta(\widehat{w}_n) - \theta(w_0)\|_X^2$  using a weighted sum of noise variables.

Lemma 1. 
$$n \cdot \|\theta(\hat{w}_n) - \theta(w_0)\|_n^2 \leq 2\sum_{i=1}^n \varepsilon_i \langle x_i, \theta(\hat{w}_n) - \theta(w_0) \rangle.$$

*Proof.* By definition of  $\hat{w}_n$ , we have  $\sum_{i=1}^n (y_i - \langle x_i, \theta(\hat{w}_n) \rangle)^2 \leq \sum_{i=1}^n (y_i - \langle x_i, \theta(w_0) \rangle)^2$ . Plugging in  $y_i = \langle x_i, \theta(w_0) \rangle + \varepsilon_i$ , breaking up the squares and cancelling out the common  $\sum_{i=1}^n \varepsilon_i^2$  terms on both sides of the inequality, we have  $\sum_{i=1}^n |\langle x_i, \theta(\hat{w}_n) - \theta(w_0) \rangle|^2 + 2\sum_{i=1}^n \varepsilon_i \langle x_i, \theta(\hat{w}_n) - \theta(w_0) \rangle \leq 0$ . Re-arranging the terms we proved the lemma.

We next adopt a localized empirical process approach (van de Geer, 2000) to upper bound  $\sum_{i=1}^{n} \varepsilon_i \langle x_i, \theta(\hat{w}_n) - \theta(w_0) \rangle$ . Define

$$\Theta_{X,F_1} := \{\theta(w) - \theta(\widetilde{w}) : w, \widetilde{w} \in \mathbb{R}^m, \|\theta(w) - \theta(\widetilde{w})\|_X \leq 1\}.$$
(16)

Re-scaling  $(w_0, \hat{w}_n) \mapsto (w_0, \hat{w}_n) / \|\theta(w) - \theta(\tilde{w})\|_X$  and  $\varepsilon_i \mapsto \varepsilon_i / \sigma$ , we have

$$\frac{1}{n}\sum_{i=1}^{n}\varepsilon_{i}\langle x_{i},\theta(\widehat{w}_{n})-\theta(w_{0})\rangle \leqslant \sigma \|\theta(\widehat{w}_{n})-\theta(w_{0})\|_{X} \cdot \sup_{\phi\in\Theta_{X,F_{1}}}\frac{1}{n}\sum_{i=1}^{n}\widetilde{\varepsilon}_{i}\langle x_{i},\phi\rangle.$$
(17)

For  $\phi \in \mathbb{R}^d$  denote  $\mathbb{G}_n^X(\phi) := \sum_{i=1}^n \tilde{\varepsilon}_i \langle x_i, \phi \rangle / \sqrt{n}$  as a random variable with randomness induced by the noise variables  $\{\tilde{\varepsilon}_i\}_{i=1}^n$ . Because  $\{\tilde{\varepsilon}_i\}_{i=1}^n$  are i.i.d. centered sub-Gaussian random variables with parameter 1, it is easy to verify that for any  $\phi, \phi' \in \mathbb{R}^d, \mathbb{G}_n^X(\phi) - \mathbb{G}_n^X(\phi')$  is a centered sub-Gaussian random variable with sub-Gaussian parameter  $\gamma^2 \leq \sigma^2 \|\phi - \phi'\|_X^2$ . Using Dudley's entropic integral (Dudley, 1967), we have

$$\mathbb{E} \sup_{\phi \in \Theta_{X,F_1}} \mathbb{G}_n^X(\phi) \lesssim \int_0^\infty \sqrt{\log N(\epsilon; \Theta_{X,F_1}, \|\cdot\|_X)} \mathrm{d}\epsilon,$$
(18)

Combining Eq. (18) with Lemma 1, we immediately have

$$\mathbb{E}\|\theta(\widehat{w}_n) - \theta(w_0)\|_X \lesssim \sqrt{\frac{1}{n}} \cdot \int_0^\infty \sqrt{\log N(\epsilon; \Theta_{X, F_1}, \|\cdot\|_X)} \mathrm{d}\epsilon.$$
(19)

In the rest of the proof we upper bound the integration of covering numbers in Eq. (19). We first consider a "population" version of the localized set  $\Theta_{X,F_1}$ :

$$\Theta_{\mu,F_1} := \{\theta(w) - \theta(\widetilde{w}) : w, \widetilde{w} \in \mathbb{R}^m, \|\theta(w) - \theta(\widetilde{w})\|_{\mu} \le 1\}$$
(20)

and upper bounds the covering number  $N(\epsilon; \Theta_{\mu,F_1}, \|\cdot\|_X)$ . We shall discuss how such an upper bound can be converted into a bound on  $N(\epsilon; \Theta_{X,F_1}, \|\cdot\|_X)$  later this section.

We first state two technical lemmas.

**Lemma 2.** For any  $\phi = \theta(w) - \theta(\tilde{w}) \in \Theta_{\mu,F_1}$  it holds that  $||w - \tilde{w}||_2 \leq \kappa^{-1}J^2$ , where J = m/s.

*Proof.* By definition, for any  $\phi = \theta(w) - \theta(\widetilde{w}) \in \Theta_{\mu,F_1}$  it holds that  $\mathbb{E}_{\mu} |\langle \theta(w) - \theta(\widetilde{w}), x \rangle|^2$ . The non-degeneracy condition (A3) of  $\mu$  then implies

$$\|\theta(w) - \theta(\widetilde{w})\|_{2}^{2} \leqslant \kappa^{-2} \mathbb{E}_{\mu} |\langle \theta(w) - \theta(\widetilde{w}), x \rangle|^{2} \leqslant \kappa^{-2}.$$
 (21)

Let  $Q_s^0 w, \ldots, Q_s^{J-1} w$  be J = m/s segments of w, each of length s. Let also  $Q_s^0 \theta, \ldots, Q_s^{r-1} \theta$  be r segments of  $\theta \in \mathbb{R}^d$ , each of length s too. Then

$$\mathsf{Q}_{s}^{\ell}[\theta(w) - \theta(\widetilde{w})] = \sum_{\ell'=0}^{\min(J-1,\ell)} \mathsf{Q}_{s}^{\ell}(w - \widetilde{w}).$$
(22)

Because  $\|\theta(w) - \theta(\widetilde{w})\|_2 \leq \kappa^{-1}$ , it holds that  $\|\mathbf{Q}_s^{\ell}[\theta(w) - \theta(\widetilde{w})]\|_2 \leq \|\theta(w) - \theta(\widetilde{w})\|_2 \leq \kappa^{-1}$  for all  $\ell \in \{0, 1, \dots, r-1\}$  because  $\mathbf{Q}_s^{\ell}[\theta(w) - \theta(\widetilde{w})]$  partitions  $\theta(w) - \theta(\widetilde{w})$  into disjoint segments.

Since  $Q_s^0[\theta(w) - \theta(\widetilde{w})] = Q_s^0(w - \widetilde{w})$ , we know that  $\|Q_s^0(w - \widetilde{w})\|_2 = \|Q_s^0[\theta(w) - \theta(\widetilde{w})]\|_2 \leq \|\theta(w) - \theta(\widetilde{w})\|_2$ . Similarly, because  $Q_s^1[\theta(w) - \theta(\widetilde{w})] = Q_s^0(w - \widetilde{w}) + Q_s^1(w - \widetilde{w})$ , we have  $\|Q_s^1(w - \widetilde{w})\|_2 \leq \|Q_s^1[\theta(w) - \theta(\widetilde{w})]\|_2 + \|Q_s^0(w - \widetilde{w})\|_2 \leq 2\|\theta(w) - \theta(\widetilde{w})\|_2$ . Continuing this argument we have  $\|Q_s^\ell(w - \widetilde{w})\|_2 \leq (\ell + 1)\|\theta(w) - \theta(\widetilde{w})\|_2$ . Subsequently,

$$\|w - \widetilde{w}\|_{2} \leq \sum_{\ell=0}^{J-1} \|\mathsf{Q}_{s}^{\ell}(w - \widetilde{w})\|_{2} \leq J^{2} \cdot \|\theta(w) - \theta(\widetilde{w})\|_{2} \leq \kappa^{-1} J^{2}.$$
 (23)

**Lemma 3.** Fix arbitrary  $\delta \in (0, 1/2)$ . With probability  $1 - \delta$  over the random draws of  $x_1, \ldots, x_n \sim \mu$ , for any  $\phi = \theta(w) - \theta(\widetilde{w}), \phi' = \theta(w') - \theta(\widetilde{w}') \in \mathbb{R}^m$  it holds that  $\|\phi - \phi'\|_X^2 \lesssim \nu^2 \sqrt{d^3 \log(1/\delta)} \cdot \|w - \widetilde{w} - w' + \widetilde{w}'\|_2^2$ .

Proof. By definition we have that

$$\|\phi - \phi'\|_X^2 = \frac{1}{n} \sum_{i=1}^n \left| \langle x_i, \phi - \phi' \rangle \right|^2 \le \lambda_{\max}(\hat{\Sigma}_n) \|\phi - \phi'\|_2^2,$$
(24)

where  $\lambda_{\max}(\hat{\Sigma}_n)$  is the largest eigenvalue of sample covariance  $\hat{\Sigma}_n = \frac{1}{n} \sum_{i=1}^n x_i x_i^{\top}$ .

Let  $\Sigma_0 := \mathbb{E}_{\mu} x x^{\top}$  denote the population covariance under  $\mu$ . Because  $\mu$  is sub-Gaussian with parameter  $\nu^2$  (A2), by standard concentration inequality of sub-Gaussian sample covariances (e.g., (Vershynin, 2012)) we have with probability  $1 - \delta$  that

$$\|\widehat{\Sigma}_n - \Sigma_0\|_{\text{op}} \lesssim \nu^2 \sqrt{d\log(1/\delta)/n}.$$
(25)

Note that  $\|\Sigma_0\|_{op}$  must also be upper bounded by  $\nu^2$  because of the sub-Gaussianity of  $\mu$ . Therefore, with probability  $1 - \delta$ 

$$\|\phi - \phi'\|_X^2 \lesssim \nu^2 \left(1 + \sqrt{\frac{d\log(1/\delta)}{n}}\right) \|\phi - \phi'\|_2^2 \lesssim \nu^2 \sqrt{d\log(1/\delta)} \cdot \|\phi - \phi'\|_2^2.$$
(26)

Finally, recall the definition that  $\phi - \phi' = \theta(w) - \theta(\widetilde{w}) - \theta(w') + \theta(\widetilde{w}') = \sum_{\ell=0}^{r-1} \mathsf{S}_s^\ell(w - \widetilde{w} - w' + \widetilde{w}')$ , implying that  $\|\phi - \phi'\|_2 \leq r \cdot \|w - \widetilde{w} - w' + \widetilde{w}'\|_2 \leq d\|w - \widetilde{w} - w' + \widetilde{w}'\|_2$ . The lemma is thus proved.

We are now ready to state and prove our key covering number lemma of  $\Theta_{\mu,F_1}$  in  $\|\cdot\|_X$ .

**Lemma 4.** With probability  $1 - \delta$  over the random draw of  $x_1, \ldots, x_n \sim \mu$ , it holds for all  $\epsilon > 0$  that  $\log N(\epsilon; \Theta_{\mu, F_1}, d_{\mu}(\cdot, \cdot)) \leq m \log(\nu d \log(\delta^{-1})/\epsilon \kappa)$ .

*Proof.* For any R > 0 denote  $\mathbb{B}_m(R) := \{z \in \mathbb{R}^m : ||z||_2 \leq R\}$  as the centered *m*-dimensional Euclidean ball of radius *R*. By Lemma 2, we know that  $\{w - \tilde{w} : \theta(w) - \theta(\tilde{w}) \in \Theta_{\mu,F_1}\} \subseteq \mathbb{B}_m(\kappa^{-1}J^2)$ .

Let  $\mathbf{H} \subseteq \mathbb{R}^m$  be a finite covering set of  $\mathbb{B}_m(\kappa^{-1}J^2)$  in  $\|\cdot\|_2$  up to a difference precision parameter  $\epsilon' > 0$  to be specified later, meaning that  $\sup_{\Delta w \in \mathbb{B}_m(\kappa^{-1}J^2)} \min_{\Delta w' \in \mathbf{H}} \|\Delta w - \Delta w'\|_2 \leq \epsilon'$ . Again using the standard covering number of *m*-dimensional unit balls (e.g., (van de Geer, 2000)), the size of  $|\mathbf{H}|$  can be upper bounded by  $\log |\mathbf{H}| \leq m \log(J^2/\kappa\epsilon')$ .

Let  $\Phi(\mathbf{H}) := \{ \phi' = \theta(w') - \theta(\widetilde{w}') : w' - \widetilde{w}' \in \mathbf{H} \}$  be the induced *d*-dimensional parameter sets by **H**. Clearly  $\log |\Phi(\mathbf{H})| \leq \log |\mathbf{H}|$ . On the other hand, by Lemma 3 and the fact that  $\{w - \widetilde{w} : \theta(w) - \theta(\widetilde{w}) \in \Theta_{\mu,F_1}\} \subseteq \mathbb{B}_m(\kappa^{-1}J^2)$ , we have  $\sup_{\phi \in \Theta_{\mu,F_1}} \min_{\phi' \in \Phi(\mathbf{H})} \|\phi - \phi'\|_X \leq \nu (d^3 \log(1/\delta))^{1/4} \cdot \epsilon'$ . Putting  $\epsilon' \approx \epsilon/\nu (d^3 \log(1/\delta))^{1/4}$  we proved the lemma.

Finally, we show how an upper bound on  $N(\epsilon; \Theta_{\mu,F_1}, \|\cdot\|_X)$  can be turned into an upper bound on  $N(\epsilon'; \Theta_{X,F_1}, \|\cdot\|_X)$  for a potentially different precision parameter  $\epsilon'$ . This is done by considering the following "restricted eigenvalue" (Bickel et al., 2009) type conditions.

**Lemma 5.** Fix arbitrary  $\delta \in (0, 1/2)$ . If  $\nu \sqrt{\log(n/\delta)} \ge \kappa$  and  $n \gtrsim \kappa^{-2}\nu^2 m \log(\kappa^{-1}\nu d \log \delta^{-1}) \log(n\delta^{-1})$ , then with probability  $1 - \delta$  we have that  $\|\phi\|_X^2 \ge 1/2 \|\phi\|_{\mu}^2$  uniformly for all  $\phi \in \Theta_{\mu, F_1}$ .

*Proof.* Because both  $\|\cdot\|_{\mu}$  and  $\|\cdot\|_X$  are linear (i.e.,  $\|a\phi\| = a\|\phi\|$  for all  $a \in \mathbb{R}$ ), it suffices to consider  $\phi \in \Theta_{\mu,F_1}$  with  $\|\phi\|_{\mu} = 1$  only.

We first consider the case of fixed  $\phi \in \Theta_{\mu,F_1}$ . Because  $\|\phi\|_{\mu}^2 = 1$ , we have  $\|\phi\|_2^2 \leq \kappa^{-2}$  thanks to (A3). In addition, because  $x_1, \ldots, x_n \sim \mu$  are independent sub-Gaussian random vectors with sub-Gaussian parameter  $\nu^2$ , we have with probability  $1 - 0.1\delta$  that  $\max_i \|x_i\|_2 \leq \nu \sqrt{\log(n/\delta)}$ . Subsequently,  $\max_i |\langle \phi, x_i \rangle|^2 \leq \kappa^{-1}\nu \sqrt{\log(n/\delta)}$ . Conditioned on this event, using Hoeffding's concentration inequality (Hoeffding, 1963) we have with probability  $1 - \delta'$  that

$$\left| \|\phi\|_X^2 - \|\phi\|_{\mu}^2 \right| = \left| \frac{1}{n} \sum_{i=1}^n |\langle x_i, \phi \rangle|^2 - \mathbb{E}_{\mu} |\langle x, \phi \rangle|^2 \right| \lesssim \kappa^{-1} \nu \sqrt{\log(n/\delta)} \cdot \sqrt{\frac{\log(1/\delta')}{n}}.$$
(27)

Next consider a finite covering set  $\mathbf{H}(\epsilon')$  of  $\Theta_{\mu,F_1}$  in distance metric  $\|\cdot\|_X$  up to a precision parameter  $\epsilon' > 0$  to be specified later; that is,  $\sup_{\phi \in \Theta_{\mu,F_1}} \min_{\phi' \in \mathbf{H}(\epsilon')} \|\phi - \phi'\|_X \leq \epsilon'$ . Lemma 4 guarantees the existence of such a covering set with size  $\log |\mathbf{H}(\epsilon')| \leq m \log(\nu d \log(\delta^{-1})/\epsilon' \kappa)$  with probability  $1 - 0.1\delta$ . On the other hand,

$$\begin{split} \left| \|\phi\|_{X}^{2} - \|\phi'\|_{X}^{2} \right| &= \left| \frac{1}{n} \sum_{i=1}^{n} \langle x_{i}, \phi \rangle^{2} - \langle x_{i}, \phi' \rangle^{2} \right| = \left| \frac{1}{n} \sum_{i=1}^{n} \langle x_{i}, \phi + \phi' \rangle \langle x_{i}, \phi - \phi' \rangle \right| \\ &\leq \sqrt{\frac{1}{n} \sum_{i=1}^{n} |\langle x_{i}, \phi + \phi' \rangle|^{2}} \sqrt{\frac{1}{n} \sum_{i=1}^{n} |\langle x_{i}, \phi - \phi' \rangle|^{2}} \\ &\leq \|\phi + \phi'\|_{X} \cdot \|\phi - \phi'\|_{X} \leq (2\|\phi\|_{X} + \|\phi - \phi'\|_{X}) \cdot \|\phi - \phi'\|_{X}. \end{split}$$

Because  $\|\phi\|_{\mu} = 1$ , we have  $\|\phi\|_X \leq \max_i \|x_i\|_2 \cdot \|\phi\|_2 \leq \kappa^{-1} \nu \sqrt{\log(n/\delta)}$  with probability  $1 - 0.1\delta$  uniformly for all  $\phi \in \Theta_{\mu, F_1}$ . Subsequently,

$$\left| \|\phi\|_X^2 - \|\phi'\|_X^2 \right| \lesssim (\kappa^{-1}\nu\sqrt{\log(n/\delta)} + \epsilon')\epsilon'.$$
(28)

Setting  $\epsilon' \simeq O(1), \delta' = 0.1\delta/|\mathbf{H}(\epsilon')|$  and combining Eqs. (27,28) we proved the desired lemma. **Corollary 1.** Under the same conditions in Lemma 4,  $N(\epsilon; \Theta_{X,F_1}, \|\cdot\|_X) \leq N(\epsilon/2; \Theta_{\mu,F_1}, \|\cdot\|_X)$ .

Combining Lemmas 1, 4 and Corollary 1 we have with probability  $1 - \delta$  that

$$\mathbb{E}\|\theta(\hat{w}_n) - \theta(w_0)\|_X \lesssim \sqrt{\frac{\sigma^2 m \log(\kappa^{-1}\nu d \log(\delta^{-1}))}{n}}.$$
(29)

Invoking Lemma 5 again we can upper bound  $\|\theta(\hat{w}_n) - \theta(w_0)\|_{\mu} = \operatorname{err}_{\mu}(\hat{w}_n, w_0; F).$ 

# **B** Proof of Theorem 2 (lower bound, average pooling)

We use the standard *Fano's inequality* (Yu, 1997; Tsybakov, 2009) to prove the minimax lower bound in Theorem 2. Below we state a commonly used variant of Fano's inequality from (Tsybakov, 2009), also known as the *Tsybakov's master theorem*:

**Lemma 6.** Let  $\mathcal{W} = (w_0, w_1, \dots, w_M)$  be a finite collection of parameters and let  $P_j$  be the distribution induced by parameter  $w_j$ , for  $j \in \{0, \dots, M\}$ . Let also  $D : \mathcal{W} \times \mathcal{W} \to \mathbb{R}^+$  be a semi-distance. Suppose the following conditions hold:

- 1.  $D(w_j, w_k) \ge 2\rho > 0$  for all  $j, k \in \{0, ..., M\}$ ;
- 2.  $P_j \ll P_0$  for every  $j \in \{1, ..., M\}; {}^4$
- 3.  $\frac{1}{M} \sum_{i=1}^{M} \operatorname{KL}(P_{j} \| P_{0}) \leq \gamma \log M;$

then the following bound holds:

$$\inf_{\hat{w}} \sup_{w_j \in \mathcal{W}} \Pr_j \left[ D(\hat{w}, w_j) \ge \rho \right] \ge \frac{\sqrt{M}}{1 + \sqrt{M}} \left( 1 - 2\gamma - 2\sqrt{\frac{\gamma}{\log M}} \right).$$
(30)

Recall that J = m/s is the number of strides in each filter, which is assumed to be an integer. We consider two cases separately.

The non-overlapping case J = 1. Construct subset  $\mathcal{Z} = \{z_1, \ldots, z_m\} \in \{-1, 1\}^m$  such that

- 1. For every  $z \in \mathbb{Z}$ ,  $\sum_{i=1}^{n} z_i = 0$ ;
- 2. For every distinct pairs of  $z, z' \in \mathbb{Z}$ ,  $\Delta_H(z, z') = \sum_{i=1}^m \mathbf{1}\{z_i \neq z'_i\} \ge d/16$ .

Using classical constructions of separable constant-weight codes (e.g., (Wang & Singh, 2016, Lemma 9), (Graham & Sloane, 1980, Theorem 7)), such a subset  $\mathcal{Z}$  exists with size  $\log |\mathcal{Z}| \gtrsim m$ .

Construct  $\mathcal{W} = \{w_0, w_1, \dots, w_M\} \subseteq \mathbb{R}^m$  as  $w_0 = 0$  and  $w_j = \delta z_j$  for  $j \in \{1, \dots, M\}$  and some  $\delta > 0$  to be specified later. Recall that  $\theta(w) = \sum_{\ell=0}^{r-1} \mathsf{S}_s^\ell w$ . Note also that  $x_1, \dots, x_n \sim \mathcal{N}(0, I)$  and  $\varepsilon_1, \dots, \varepsilon_n \sim \mathcal{N}(0, 1)$ . The conditions in Lemma 6 can be verified below:

- 1. For every  $w, w' \in \mathcal{W}, D(w, w') := \|\theta(w) \theta(w')\|_2 \ge \sqrt{r\delta} \cdot d/12 \ge \delta\sqrt{rd};$
- 2. For every  $w_j \in \{w_1, \ldots, w_M\}$ , we have  $P_j \ll P_0$  and furthermore  $\operatorname{KL}(P_j || P_0) = \mathbb{E} \sum_{i=1}^n |\langle x_i, \theta(w_j) \theta(w_0) \rangle|^2 / 2\sigma^2 = n ||\theta(w_j) \theta(w_0)||_2^2 / \sigma^2 \lesssim nr\delta^2 / \sigma^2$ .

Setting  $\delta \approx \sqrt{\sigma^2 m/nr}$  and invoking Lemma 6 we have

$$\inf_{\hat{w}} \sup_{w_j \in \mathcal{W}} \Pr_j \left[ \| \hat{w} - w_j \|_2 \ge c_0 \sqrt{\frac{\sigma^2 m}{n}} \right] \ge \frac{1}{4}.$$
(31)

Theorem 2 is then proved by applying the Markov's inequality and noting that  $\operatorname{err}_{\mu}(\hat{w}, w) = \sqrt{\mathbb{E}_{\mu}|\langle x, \hat{w} - w \rangle|^2} = \|\hat{w} - w\|_2.$ 

 $<sup>{}^{4}</sup>P \ll Q$  means that the support of P is contained in the support of Q.

The overlapping case J > 1. Construct subset  $\mathcal{Z} = \{z_1, \ldots, z_m\} \in \{-1, 1\}^m$  such that

- 1. For every  $z \in \mathcal{Z}$ ,  $\sum_{i=1}^{n} z_i = 0$  and  $Q_s^{J-1} z = 0$ ;
- 2. For every distinct pairs of  $z, z' \in \mathbb{Z}$ ,  $\Delta_H(z, z') = \sum_{i=1}^m \mathbf{1}\{z_i \neq z'_i\} \ge d/16$ .

Using again the construction of separable constant-weight codes, such a subset Z exists with size  $\log |Z| \gtrsim m - s \ge m/2$ .

For every  $z \in \mathbb{Z} \subseteq \mathbb{R}^m$ , construct  $w(z) \in \mathbb{R}^m$  as follows:

$$\mathbf{Q}_{s}^{\ell}w := \mathbf{Q}_{s}^{\ell}z - \sum_{\ell'=0}^{\ell-1} \mathbf{Q}_{s}^{\ell'}z \qquad \ell = 0, 1, \dots, J-1,$$
(32)

It is then easy to verify that  $\theta(w(z)) = \sum_{\ell=0}^{r-1} \mathsf{S}_s^\ell w(z) = (z_1, \dots, z_m, 0, \dots, 0).$ 

Construct  $\mathcal{W} = \{w_0, w_1, \dots, w_M\} \subseteq \mathbb{R}^m$  as  $w_0 = 0$  and  $w_j = \delta' z_j$  for  $j \in \{1, \dots, M\}$  and some  $\delta' > 0$  to be specified later. The analysis in the non-overlapping case remains valid:

- 1. For every  $w, w' \in \mathcal{W}$ ,  $D(w, w') := \|\theta(w) \theta(w')\|_2 \Rightarrow \delta' \cdot d/12 \gtrsim \delta \sqrt{r} d;$
- 2. For every  $w_j \in \{w_1, \ldots, w_M\}$ , we have  $P_j \ll P_0$  and furthermore  $\operatorname{KL}(P_j || P_0) = \mathbb{E} \sum_{i=1}^n |\langle x_i, \theta(w_j) \theta(w_0) \rangle|^2 / 2\sigma^2 = n ||\theta(w_j) \theta(w_0)||_2^2 / \sigma^2 \leq n\delta^2 / \sigma^2$ .

Setting  $\delta' = \sqrt{\sigma^2 m/n}$  we complete the proof.

# **C Proof of Theorem 3 (upper bound, prediction layers)**

We use a similar framework as in the proof of Theorem 1 to prove Theorem 3. For the convolutional network with prediction layers, the parameterization  $\theta(w, a) \in \mathbb{R}^d$  takes the form of

$$\theta(w,a) = \sum_{\ell=0}^{r-1} a_{\ell} \mathsf{S}_{s}^{\ell} w.$$
(33)

Deifne

$$\Theta_{X,F_2} := \{\theta(w,a) - \theta(\widetilde{w},\widetilde{a}) : w, \widetilde{w} \in \mathbb{R}^m, a, \widetilde{a} \in \mathbb{R}^r, \|\theta(w,a) - \theta(\widetilde{w},\widetilde{a})\|_X \le 1\}.$$
(34)

Using the same basic inequality and Dudley's entropic integral as in the proof of Theorem 1, we have

$$\mathbb{E}\|\theta(\hat{w}_n, \hat{a}_n) - \theta(w_0, a_0)\|_X \lesssim \sqrt{\frac{1}{n}} \cdot \int_0^\infty \sqrt{\log N(\epsilon; \Theta_{X, F_2}, \|\cdot\|_X)} \mathrm{d}\epsilon.$$
(35)

We similarly also consider a population version of  $\Theta_{X,F_1}$ :

$$\Theta_{\mu,F_2} := \{\theta(w,a) - \theta(\widetilde{w},\widetilde{a}) : w, \widetilde{w} \in \mathbb{R}^m, a, \widetilde{a} \in \mathbb{R}^r, \|\theta(w,a) - \theta(\widetilde{w},\widetilde{a})\|_{\mu} \le 1\}.$$
 (36)

The following lemma upper bounds the covering number of  $\Theta_{\mu,F_2}$  with respect to  $\|\cdot\|_X$ .

**Lemma 7.** With probability  $1 - \delta$  over the random draw of  $x_1, \ldots, x_n \sim \mu$ , it holds for all  $\epsilon > 0$  that  $\log N(\epsilon; \Theta_{\mu, F_2}, \|\cdot\|_X) \leq (rJ + m) \log(\kappa^{-1}\nu d \log \delta^{-1})$ .

*Proof.* Consider any  $\phi = \theta(w, a) - \theta(\tilde{w}, \tilde{a}) \in \Theta_{\mu, F_2}$ . By definition, we know that  $\|\phi\|_{\mu} \leq 1$ , and therefore  $\|\phi\|_2 \leq \kappa^{-1}$  thanks to the non-degeneracy condition (A3).

Let  $Q_s^0 \phi, \ldots, Q_s^{r-1} \phi$  be the *r* disjoint segments of  $\phi \in \mathbb{R}^d$ , each of length *s*. Let also  $Q_s^0 w, \ldots, Q_s^{J-1} w$  be the *J* disjoint segments of  $w \in \mathbb{R}^m$  each of length *s*. Denote for convenience that  $a_t = a_t \mod r$  and  $\tilde{a}_t = \tilde{a}_t \mod r$ , which extends the subscripts of *a* and  $\tilde{a}$  to  $\mathbb{Z}$ . Then

$$\mathsf{Q}_{s}^{\ell}\phi = \sum_{j=0}^{J-1} a_{\ell+j}\mathsf{Q}_{s}^{j}w - \widetilde{a}_{\ell+j}\mathsf{Q}_{s}^{j}\widetilde{w}, \qquad \ell = 0, 1, \dots, r-1.$$
(37)

Subsequently, for every  $\ell$  we have  $Q_s^{\ell}\phi \in \operatorname{span}\{Q_s^jw, Q_s^j\widetilde{w}\}_{j=0}^{J-1}$ , a linear subspace in  $\mathbb{R}^s$  of dimension at most 2*J*. Note also that  $\|Q_s^{\ell}\phi\|_2 \leq \|\phi\|_2 \leq \kappa^{-1}$ . We then have (recall that  $\mathbb{B}_s(R)$  denotes the centered *m*-dimensional Euclidean ball of radius *R*)

$$\Theta_{\mu,F_2} \subseteq \{\phi \in \mathbb{R}^m : \exists \mathcal{S} \subseteq \mathbb{R}^s, \dim(S) \leqslant J \ s.t. \ \mathsf{Q}_s^\ell \phi \in \mathcal{S} \cap \mathbb{B}_s(\kappa^{-1})\} =: \widetilde{\Theta}.$$
(38)

Our construction of covering sets of  $\widetilde{\Theta}$  (and therefore also  $\Theta_{\mu,F_2}$ ) can be divided into two steps. As a first step, we construct covering set  $\mathbf{S} = \{\mathcal{S}_1^*, \ldots, \mathcal{S}_N^*\}$  such that each  $\mathcal{S}_k^*, \ell \in [N]$  is linear subspace in  $\mathbb{R}^s$  of dimension at most 2J, and furthermore for any linear subspace  $\mathcal{S} \subseteq \mathbb{R}^s$ ,  $\dim(\mathcal{S}) \leq 2J$ ,

$$\min_{\mathcal{S}_{k}^{*} \in \mathbf{S}} \sup_{z \in \mathcal{S} \cap \mathbb{B}_{s}(1)} \inf_{z' \in \mathcal{S}_{k}^{*} \cap \mathbb{B}_{s}(1)} \|z - z'\|_{2} \leqslant \epsilon',$$
(39)

where  $\epsilon' > 0$  is an error tolerance parameter to be specified later. The following proposition gives an upper bound on the size of such coverings, which is proved later.

**Proposition 1.** There exists **S** satisfying Eq. (39); furthermore,  $\log |\mathbf{S}| \leq m \log(d/\epsilon')$ .

The next step is to construct, for each  $S_k^* \in \mathbf{S}$ , a covering  $\mathbf{H}(S_k^*) = \{u_1^*, \dots, u_T^*\} \subseteq S_k^* \cap \mathbb{B}_s(1)$  satisfying

$$\sup_{\substack{u \in \mathcal{S}_k^* \cap \mathbb{B}_s(1) \ u_t^* \in \mathbf{H}(\mathcal{S}_k^*)}} \|u - u_t^*\|_2 \le \epsilon'',\tag{40}$$

where  $\epsilon'' > 0$  is another error tolerance parameter to be specified later. The following proposition gives an upper bound on the size of such coverings. Its proof is also given later.

**Proposition 2.** There exists  $\mathbf{H}(\mathcal{S}_k^*)$  satisfying Eq. (40); furthermore,  $\log |\mathbf{H}(\mathcal{S}_k^*)| \leq J \log(1/\epsilon'')$ .

We now construct our final covering set as follows:

$$\mathbf{C} := \bigcup_{\mathcal{S}_{k}^{*} \in \mathbf{S}} \{ \phi \in \mathbb{R}^{d} : \mathsf{Q}_{s}^{0} \phi, \dots, \mathsf{Q}_{s}^{r-1} \phi \in \mathbf{H}(\mathcal{S}_{k}^{*}) \}.$$
(41)

For any  $\phi \in \widetilde{\Theta}$  corresponding to linear subspace S in  $\mathbb{R}^s$ , dim $(S) \leq 2J$ , one first finds  $S_k^* \in \mathbf{S}$  that best approximates S in the sense of Eq. (39). Then for each  $Q_s^\ell \phi$ ,  $\ell \in \{0, \ldots, r-1\}$ , one can find  $u_\ell^* \in \mathbf{H}(S_k^*)$  that minimizes  $\|Q_s^\ell \widetilde{\phi} - \widetilde{u}_\ell^*\|_2$  where  $\widetilde{u}_\ell^* = \|Q_s^\ell \widetilde{\phi}\|_2 \widetilde{u}_\ell^*$  and  $Q_s^\ell \widetilde{\phi} \in S_\ell^*$  minimizes  $\|Q_s^\ell (\phi - \widetilde{\phi})\|_2$ . Define  $\phi^* \in \mathbb{R}^d$  as  $Q_s^\ell \phi^* = \widetilde{u}_\ell^*$ . Then

$$\|\phi - \phi^*\|_2 \leq \sum_{\ell=0}^{r-1} \|\mathbf{Q}_s^{\ell}(\phi - \widetilde{\phi})\|_2 + \|\mathbf{Q}_s^{\ell}\widetilde{\phi} - \widetilde{u}_\ell^*\|_2 \leq r\kappa^{-1}(\epsilon' + \epsilon''),$$
(42)

where the last inequality holds by a scaling argument and the fact that  $\|Q_s^\ell \phi\|_2 \leq \kappa^{-1}$ . Finally, because  $x_1, \ldots, x_n \sim \mu$  are independent sub-Gaussian random vectors, using Eq. (25) we have with probability  $1 - \delta$  that

$$\|\phi - \phi^*\|_X \lesssim \nu^2 \sqrt{d \log(1/\delta)} \cdot \|\phi - \phi^*\|_2 \leqslant \kappa^{-1} \nu^2 \sqrt{d^3 \log(1/\delta)} \cdot (\epsilon' + \epsilon'').$$
(43)

Setting  $\epsilon' = \epsilon'' \approx \epsilon/(\kappa^{-1}\nu^2 \sqrt{d^3 \log(1/\delta)})$  we proved that **C** is a valid covering set of  $\widetilde{\Theta}$  (and also  $\Theta_{\mu,F_2}$ ) in  $\|\cdot\|_X$  up to precision  $\epsilon$ .

Finally we count the number of elements in C. By Propositions 1 and 2, we have  $\log |\mathbf{S}| \leq J \log(\kappa^{-1}\nu d \log \delta^{-1})$  and  $\log |\mathbf{H}(\mathcal{S}_{\ell}^*)| \leq J \log(\kappa^{-1}\nu d \log \delta^{-1})$ . By construction of C, we have  $|\mathbf{C}| \leq |\mathbf{S}| \times \max_{\mathcal{S}_{k}^{*} \in \mathbf{S}} |\mathbf{H}(\mathcal{S}_{k}^{*})|^{r}$ . Subsequently,  $\log |\mathbf{C}| \leq (rJ + m) \log(\kappa^{-1}\nu d \log \delta^{-1})$ .  $\Box$ 

*Proof of Proposition 1.* If  $2J \ge s$  then the proposition clearly holds. So we shall only prove the proposition in cases where  $2J \le s$ .

Let  $\mathcal{U}, \mathcal{V}$  be two linear subspace of  $\mathbb{R}^s$  of dimension at most 2J. Let  $U, V \in \mathbb{R}^{s \times 2J}$  be the corresponding orthonormal basis of  $\mathcal{U}$  and  $\mathcal{V}$ , with orthogonal columns. Any  $u \in \mathcal{U} \cap \mathbb{B}_s(1)$  can then be written as  $u = U\alpha$  with  $\|\alpha\|_2 = 1$ . Consider  $v := V\alpha$ . It is easy to verify that  $v \in \mathcal{V} \cap \mathbb{B}_s(1)$ . In addition,  $\|u - v\|_2 = \|(U - V)\alpha\|_2 \leq \|U - V\|_{\text{op}} \leq \|U - V\|_F$ . Subsequently, a covering of  $\{U \in \mathbb{R}^{2J \times s} : \|U\|_F \leq \sqrt{2J}\|U\|_{\text{op}} = \sqrt{2J}\}$  in  $\|\cdot\|_F$  up to precision  $\epsilon'$  implies a covering in the sense of Proposition 1. By viewing U as a  $(2J \times s)$ -dimensional vector in the Euclidean space, it is easy to see that such a cover exists with size  $\log N \leq (2Js) \log(2J/\epsilon') \leq Js \log(d/\epsilon') = m \log(d/\epsilon')$ .  $\Box$ 

*Proof of Proposition 2.* If  $2J \ge s$  then the proposition clearly holds, because one only needs to invoke classical coverings of  $\mathbb{B}_s(1)$ . So we shall only consider cases where  $2J \le s$ .

Let  $U \in \mathbb{R}^{2J \times s}$  be an orthonormal basis of  $\mathcal{S}_k^*$ . Any  $u \in \mathcal{S}_k^* \cap \mathbb{B}_s(1)$  can be written as  $u = U\alpha$ for some  $\|\alpha\|_2 \leq 1$ . For any other  $\beta \in \mathbb{R}^{2J}$ ,  $\|\beta\|_2 \leq 1$ , consider  $v = U\beta$ . It is easy to verify that  $v \in \mathcal{S}_k^* \cap \mathbb{B}_s(1)$ . Furthermore,  $\|u - v\|_2 = \|U(\alpha - \beta)\|_2 \leq \|\alpha - \beta\|_2$ . Therefore, a covering of  $\mathbb{B}_{2J}(1)$  in  $\|\cdot\|_2$  up to precision  $\epsilon''$  implies a covering in the sense of Proposition 2. On the other hand, standard covering number arguments (e.g., (van de Geer, 2000)) show that such a covering exists with size  $\log N \leq 2J \log(1/\epsilon'')$ .

Finally, we show that an upper bound on  $\log N(\epsilon; \Theta_{\mu,F_1}, \|\cdot\|_X)$  implies an upper bound on  $\log N(\epsilon'; \Theta_{X,F_1}, \|\cdot\|_X)$  for a potentially different precision parameter  $\epsilon'$ . Similar to the proof of Theorem 1, the following lemma establishes a restricted eigenvalue type condition for network  $F_2$  and the parameter spaces  $\Theta_{\mu,F_1}, \Theta_{X,F_1}$  it induces.

**Lemma 8.** Fix arbitrary  $\delta \in (0, 1/2)$ . If  $\nu \sqrt{\log(n/\delta)} \ge \kappa$  and  $n \ge \kappa^{-2}\nu^2(rJ + m)\log(\kappa^{-1}\nu d\log\delta^{-1})\log(n\delta^{-1})$ , then with probability  $1 - \delta$  we have that  $\|\phi\|_X^2 \ge 1/2\|\phi\|_{\mu}^2$  uniformly for all  $\phi \in \Theta_{\mu, F_2}$ .

The proof of Lemma 8 is identical to the proof of Lemma 5 except that a different covering number lemma is invoked; therefore we omit the proof. Combining Lemmas 7, 8 with Eq. (35) we proved Theorem 3.

# **D** Proof of Theorem 4 (lower bound, prediction layers)

Because  $r+m \leq 2 \max(r, m)$ , it suffices to prove minimax lower bounds of  $\sqrt{\sigma^2 m/n}$  and  $\sqrt{\sigma^2 r/n}$  separately.

First consider  $a_0 = (1, 0, ..., 0)$  and  $w_0 \in \mathbb{R}^m$  free to vary. Then  $F_2(x; a_0, w_0) = w_0^\top \mathsf{P}_s^0 x$  reduces to a standard linear regression problem with m covariates. It is a classical result (e.g., (Van der Vaart, 1998)) that the minimax mean-square error of learning an m-dimensional linear predictor is m/n; more specifically,

$$\inf_{\hat{w}_n} \sup_{w_0 \in \mathbb{R}^m} \mathbb{E} \| \hat{w}_n - w \|_2 \gtrsim \sqrt{\sigma^2 m/n}.$$
(44)

On the other hand, because  $a_0 = (1, 0, ..., 0)$  and  $x_1, ..., x_n \sim \mathcal{N}(0, I)$ , we have that  $\operatorname{err}^2_{\mu}(\hat{\boldsymbol{w}}_n, \boldsymbol{w}_0; F_2) = \|\hat{\boldsymbol{w}}_n - \boldsymbol{w}_0\|_2^2$ . Therefore, Eq. (44) implies a  $\sqrt{\sigma^2 m/n}$  lower bound on the minimax mean-square error  $\operatorname{Eerr}_{\mu}(\hat{\boldsymbol{w}}_n, \boldsymbol{w}_0; F_2)$ .

Next consider  $w_0 = (1, 0, ..., 0)$  and  $a \in \mathbb{R}^r$  free to vary. Denote  $\tilde{x} := (x_0, x_s, ..., x_{(r-1)s}) \in \mathbb{R}^r$ . Then  $F_2(x; a_0, w_0) = a_0^\top \tilde{x}$ . Also note that  $\tilde{x} \sim \mathcal{N}(0, I_r)$  because  $x \in \mathcal{N}(0, I_d)$ . Using the same analysis above we can establish a  $\sqrt{\sigma^2 r/n}$  lower bound on the minimax mean-square error  $\mathbb{E}\operatorname{err}_{\mu}(\hat{w}_n, w_0; F_2)$ .

Combining both cases we complete the proof of Theorem 4.