

# Robust Generalization of Quadratic Neural Networks via Function Identification

Kan Xu <sup>\*</sup>      Hamsa Bastani <sup>†</sup>      Osbert Bastani <sup>‡</sup>

## Abstract

A key challenge facing deep learning is that neural networks are often not robust to shifts in the underlying data distribution. We study this problem from the perspective of the statistical concept of parameter identification. Generalization bounds from learning theory often assume that the test distribution is close to the training distribution. In contrast, if we can identify the “true” parameters, then the model generalizes to arbitrary distribution shifts. However, neural networks typically have internal symmetries that make parameter identification impossible. We show that we can identify the function represented by a quadratic network even though we cannot identify its parameters; we extend this result to neural networks with ReLU activations. Thus, we can obtain robust generalization bounds for neural networks. We leverage this result to obtain new bounds for contextual bandits and transfer learning with quadratic neural networks. Overall, our results suggest that we can improve robustness of neural networks by designing models that can represent the true data generating process.

## 1 Introduction

Recent work has shown that neural networks are not robust to shifts in the data, including both distribution shifts (where the data comes from a new distribution independent of the neural network parameters) [25, 44] and adversarial shifts (where the shift can depend on the parameters) [43]. Accordingly, there has been interest in better understanding why neural networks fail to be robust [45, 27] and on improving robustness [24, 36, 13].

From the perspective of learning theory, there is little reason to expect neural networks to be robust, since generalization bounds typically assume that the test examples are from the same distribution as the training examples. PAC-Bayesian generalization bounds allow for a limited amount of robustness, but only if the support of the target distribution  $q$  is contained in that of the source distribution  $p$ , since it requires that the KL divergence  $D_{\text{KL}}(q \parallel p)$  is small. Yet, distribution shifts [25] often shift probability mass to inputs completely outside the source.

Instead, the reason we might expect neural networks to be robust to these shifts is that humans are robust to them; for instance, small pixel-level shifts considered in adversarial examples are typically unnoticeable to humans, yet these shifts can move the image completely off of the distribution of natural images. This fact indicates a gap in our theoretical understanding of neural networks. In particular, the key question is understanding settings under which we may expect neural networks to be robust to distribution shifts that are “large” (e.g., in terms of KL divergence).

We study a strategy for closing this gap based on the statistical concept of *identifiability* [26]. At a high level, this concept assumes that the true model belongs to the model family; then, in the

---

<sup>\*</sup>University of Pennsylvania, Department of Economics. Email: [kanxu@sas.upenn.edu](mailto:kanxu@sas.upenn.edu).

<sup>†</sup>Wharton School, Department of Operations Information and Decisions. Email: [hamsab@wharton.upenn.edu](mailto:hamsab@wharton.upenn.edu).

<sup>‡</sup>University of Pennsylvania, Department of Computer and Information Science. Email: [obastani@seas.upenn.edu](mailto:obastani@seas.upenn.edu).

limit of infinite training data, the learning algorithm can exactly recover the parameters of the true model. For instance, in linear regression, the data is generated according to the model  $y = \langle \theta^*, x \rangle + \xi$ , where  $\xi$  is  $\sigma$ -subgaussian noise. Then, under mild assumptions on the training data  $Z = (X, Y)$ , the ordinary least squares (OLS) estimator  $\hat{\theta}(Z)$  recovers the true parameter—i.e., in the limit of infinite data,  $\hat{\theta}(Z) = \theta^*$ . With finite samples, OLS satisfies high-probability convergence rates of the form

$$\|\hat{\theta}(Z) - \theta^*\|_2 \leq \epsilon. \quad (1)$$

The connection to robustness is that if (1) holds, then for *any* input  $x$  such that  $\|x\|_2 \leq x_{\max}$ , we have

$$|\langle \hat{\theta}(Z), x \rangle - \langle \theta^*, x \rangle| \leq \|\hat{\theta}(Z) - \theta^*\|_2 \|x\|_2 \leq \epsilon x_{\max}. \quad (2)$$

Thus, for *any* distribution  $q(x)$  with support on  $B_2(0, x_{\max}) = \{x \in \mathcal{X} \mid \|x\|_2 \leq x_{\max}\}$ ,  $\hat{\theta}(Z)$  obtains bounded error—i.e.,  $\mathbb{E}_{q(x)}[(\langle \hat{\theta}(Z), x \rangle - \langle \theta^*, x \rangle)^2] \leq \epsilon^2 x_{\max}^2$  with high probability.

A natural question is whether we can obtain similar kinds of parameter identification bounds for neural networks. A key complication is that neural networks parameters have symmetries that make identification impossible, since different parameters can yield the same model. Nevertheless, it may be possible to obtain bounds of the form (2)—even if we do not recover the true parameters  $\theta^*$ , we can still recover the function  $f_{\theta^*}(x)$ , which we call *function identification*.

We prove that quadratic neural networks (QNNs) satisfy function identification bounds under mild conditions. To demonstrate the utility of this result, we show how function identification can be leveraged to obtain regret guarantees for a bandit [39] where each arm is a QNN. Linear bandits fundamentally involve covariate shift since their “covariates” are arms, which are adaptively chosen through the learning process as a function of past observations; thus, existing approaches have all operated in the setting where there is a *unique* and identifiable global minimizer. Similarly, we build on recent work proving bounds on transfer learning in the setting of bounded label shift and unbounded covariate shift [9, 49]; again, we show that we can leverage function identification to easily transfer learn QNNs. Additionally, we show function identification for the subclass ReLU networks where each component is a unit vector. Finally, in Appendix B, we study implications of these results specifically for compositional generalization, which has received recent interest [4, 32].

**Related work.** Prior work has connected misspecification (i.e., the true model is in the model family) and robustness to covariate shift [40, 48]; however, having a correctly specified model is insufficient if the true parameters are not identifiable—e.g., in linear regression, if the covariance matrix  $\Sigma = \mathbb{E}_{p(x)}[xx^\top]$  is singular, then  $\theta$  is not identifiable; thus, the estimated model may not be robust. QNNs cannot be identified even if the model is correctly specified since the parameters have a continuous symmetry (i.e., orthogonal transformations).

Recent work has studied learning under adversarial examples [24, 36, 13] and corrupted training data [42, 14]. In contrast, we are interested in robustness to covariate shift; there has been recent work empirically showing that neural networks are sensitive to distribution shift [25, 44, 38, 37, 31]. Distributionally robust optimization enables training of models robust to small shifts [17], but we are interested in potentially large shifts. Unsupervised domain adaptation [10, 11] learns a model on a covariate shifted target distribution; however, they rely on unlabeled examples from the target domain, whereas we do not. There has been recent theory on robustness to adversarial perturbations—e.g., showing there may be a tradeoff between robustness and on-distribution generalization [45], and that non-robust algorithms tend to learn predictive but brittle representations compared to adversarially robust ones [27]. In contrast, we show that these tradeoffs are mitigated when the true model function can be identified despite over-parameterization. Furthermore, adversarial shifts are typically bounded (e.g., small  $\ell_\infty$  norm), whereas the shifts we consider may be large.

There has been a great deal of recent work on deep learning theory, including on QNNs and ReLU networks; however, it has largely focused on optimization [21, 28, 5, 16, 19, 41, 33, 23], and on-distribution generalization [35, 15, 28, 6, 34, 22]. In contrast, we are interested in out-of-distribution generalization.

We discuss additional related work on matrix factorization and multi-armed bandits in Appendix A, as well as a discussion of the novelty of our results.

## 2 Problem Formulation

We consider a model  $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ , with covariates  $\mathcal{X} \subseteq \mathbb{R}^d$ , labels  $\mathcal{Y} \subseteq \mathbb{R}$ , and parameters  $\theta \in \Theta \subseteq \mathbb{R}^m$ . Generalization bounds from learning theory typically have form

$$\mathbb{P}_{p(Z)}[L_p(\hat{\theta}(Z)) \leq \epsilon] \geq 1 - \delta \quad \text{where} \quad L_p(\theta) = \mathbb{E}_{p(x)}[(f_\theta(x) - f_{\theta^*}(x))^2], \quad (3)$$

where  $\epsilon, \delta \in \mathbb{R}_{>0}$ ,  $Z = \{(x_1, y_1), \dots, (x_n, y_n)\} \subseteq \mathcal{X} \times \mathcal{Y}$  with  $y_i = f_{\theta^*}(x_i) + \xi_i$  is a training set of i.i.d. observations from a distribution  $p$  (i.e.,  $p(Z) = p(x_1, y_1) \cdot \dots \cdot p(x_n, y_n)$ ),  $\xi_i$  is bounded random noise independent of  $x_i$  with  $|\xi_i| \leq \xi_{\max}$ ,  $\theta^* \in \Theta$  are the true parameters, and

$$\hat{\theta}(Z) = \arg \min_{\theta \in \Theta} \hat{L}(\theta; Z) \quad \text{where} \quad \hat{L}(\theta; Z) = \frac{1}{n} \sum_{i=1}^n (f_\theta(x_i) - y_i)^2$$

is an estimator based on the training data  $Z$ .<sup>1</sup> In particular, they assume that the training inputs  $x_i \sim p$  are i.i.d. samples from the same distribution as the test example  $x \sim p$ .

**Definition 1.** The model  $f_\theta$  and distribution  $p$  satisfy *function identification* if for any  $\epsilon, \delta \in \mathbb{R}_{>0}$ , we have  $\mathbb{P}_{p(Z)}[(f_{\hat{\theta}(Z)}(x) - f_{\theta^*}(x))^2 \leq \epsilon, \forall x \in \mathcal{X}] \geq 1 - \delta$  for  $n = |Z|$  sufficiently large.

Function identification implies generalization bounds even when the test data comes from a different distribution  $q$ . In particular, we say  $f_\theta$  *robustly generalizes* if for any  $q$  with support on  $\mathcal{X}$ , we have

$$\mathbb{P}_{p(Z)}[L_q(\hat{\theta}(Z)) \leq \epsilon] \geq 1 - \delta, \quad (4)$$

where the difference from (3) is highlighted in red. It is easy to see that function identification implies (4). Note that the true model  $f_{\theta^*}$  does not change, so there is no label shift.

## 3 Function Identification of QNNs

Traditional statistical bounds on parameter identification can provide guarantees for arbitrary covariate shift. In particular, suppose we have a bound of the form

$$\mathbb{P}_{p(Z)} \left[ \|\hat{\theta}(Z) - \theta^*\|_2 \leq \epsilon \right] \geq 1 - \delta, \quad (5)$$

and assume that the model family  $f_\theta$  is  $K$ -Lipschitz continuous in  $\theta$ ; then, we have

$$L_q(\hat{\theta}(Z)) \leq K^2 \cdot \|\hat{\theta}(Z) - \theta^*\|_2^2 \leq K^2 \epsilon^2 \quad (6)$$

with probability at least  $1 - \delta$  according to  $p(Z)$ . In particular, this bound holds for any covariate distribution  $q$ . Our goal is to extend these techniques to quadratic neural networks (QNNs), which are over-parameterized so we cannot identify the true parameters  $\theta^*$ —i.e., (5) does not hold.

<sup>1</sup>In (3), the loss  $L_p$  omits the label errors  $\xi$ ; including it would result in an additive constant to  $L_p$ . This choice ensures that the optimal parameters have zero loss—i.e.,  $L_q(\theta^*) = 0$  for any  $q$ .

### 3.1 Quadratic Neural Networks

We consider a quadratic neural network  $f_\theta$ , where  $\theta \in \mathbb{R}^{d \times k}$ , with a single hidden layer with  $k$  neurons—i.e.,  $f_\theta(x) = \sum_{j=1}^k a_j \cdot \sigma(\langle \theta_j, x \rangle)$ . We consider the over-parameterization case where  $k$  can be much larger than  $d$ . Following prior work [15], we assume that  $f_\theta$  has quadratic activations and output weights equal to one—i.e.,  $\sigma(z) = z^2$  and  $a_j = 1$  for each  $j \in [k]$ , so

$$f_\theta(x) = \sum_{j=1}^k \langle \theta_j, x \rangle^2.$$

We assume the true (training) loss is the mean-squared error  $L_p(\theta) = \mathbb{E}_{p(x)}[(f_\theta(x) - f_{\theta^*}(x))^2]$ , and we consider a model trained using an empirical estimate of this loss on the training dataset:

$$\hat{\theta}(Z) = \arg \min_{\theta \in \Theta} \hat{L}(\theta; Z) \quad \text{where} \quad \hat{L}(\theta; Z) = \frac{1}{n} \sum_{i=1}^n (f_\theta(x_i) - y_i)^2.$$

Now, our goal is to obtain a bound of the form (6); to this end, we assume the following:

**Assumption A1.**  $\|x\|_2 \leq x_{\max}$  and  $\|\theta\|_F \leq \theta_{\max}$ .

**Assumption A2.** There exists  $\alpha \in \mathbb{R}_{>0}$  such that  $\mathbb{E}_{p(x)}[(x^\top \Delta x)^2] \geq \alpha \|\Delta\|_F^2$  for any symmetric  $\Delta \in \mathbb{R}^{d \times d}$ .

Our second assumption is standard; in particular, it is closely related to the assumption in linear regression that the minimum eigenvalue of the covariance matrix is lower bounded—i.e.,  $\Sigma = \mathbb{E}_{p(x)}[xx^\top] \succ 0$ . As an example, when  $x$  is i.i.d. uniform in each component, e.g.,  $p(x) = \prod_{i=1}^d \text{Uniform}(x_i; [-1/2, 1/2])$ , then we can take  $\alpha = 1/180$ ; we give a proof in Appendix C.1.

### 3.2 Robust Generalization

Our approach leverages the fact that  $f_\theta(x) = x^\top (\theta\theta^\top)x$ ; thus,  $f_\theta$  resembles a matrix factorization model. Recent work has leveraged this connection to translate matrix factorization theory to QNNs [15]. We let  $g(\theta) = \theta\theta^\top$  and  $\tilde{f}_\phi(x) = x^\top \phi x$ , where  $\phi \in \Phi \subseteq \mathbb{R}^{d \times d}$ , in which case  $f_\theta(x) = \tilde{f}_{g(\theta)}(x)$ ; in addition, we define  $\tilde{L}_p(\phi) = \mathbb{E}_{p(x)}[(\tilde{f}_\phi(x) - \tilde{f}_{\phi^*}(x))^2]$ , where  $\phi^* = g(\theta^*)$ , and  $\hat{L}(\phi; Z) = n^{-1} \sum_{i=1}^n (\tilde{f}_\phi(x_i) - y_i)^2$ , so  $L_p(\theta) = \tilde{L}_p(g(\theta))$  and  $\hat{L}(\theta; Z) = \hat{L}(g(\theta); Z)$ . We also assume  $\|\phi\|_F \leq \phi_{\max}$ ; in general, we have  $\phi_{\max} \leq \theta_{\max}^2$  by Assumption A1.

We begin by stating several lemmas establishing the properties needed for function identification. Our first lemma says that the loss is strongly convex in  $\phi$ .

**Lemma 2.** *Under Assumption A2, the loss  $\tilde{L}_p(\phi)$  is  $2\alpha$ -strongly convex in  $\phi$ .*

We give a proof in Appendix C.2. Our next lemma says that our model family is Lipschitz in  $\phi$ .

**Lemma 3.** *Under Assumptions A1 & A2,  $\tilde{f}_\phi$  and  $\tilde{L}$  are  $K$ -Lipschitz in  $\phi$ , where  $K = 4\phi_{\max}x_{\max}^4$ .*

We give a proof in Appendix C.3. Our final lemma says that our estimate of the loss function is a uniformly good approximation of the true loss.

**Lemma 4.** *Under Assumptions A1 & A2, for any  $\delta \in \mathbb{R}_{>0}$ , we have*

$$\mathbb{P}_{p(Z)} \left[ \sup_{\theta \in \Theta} |\hat{L}(\theta; Z) - L_p(\theta) - \sigma(Z)| \leq \epsilon \right] \geq 1 - \delta,$$

where  $\sigma(Z) = n^{-1} \sum_{i=1}^n \xi_i^2$ , and letting  $\ell_{max} = 2x_{max}^2 \phi_{max}$  be an upper bound on  $|f_\theta(x) - f_{\theta^*}(x)|$ ,

$$\epsilon = \sqrt{\frac{18\ell_{max}^2(\ell_{max}^2 + \xi_{max}^2)(d^2 \max\left\{1, \log\left(1 + \frac{8\phi_{max}Kn}{\ell_{max}^2}\right)\right\} + \log\frac{2}{\delta})}{n}}. \quad (7)$$

We give a proof in Appendix C.4. Note that  $\epsilon \rightarrow 0$  as  $n \rightarrow \infty$ . Next, we prove our main result, which says that quadratic neural networks can be functionally identified.

**Theorem 5.** *Under Assumptions A1 & A2, we have*

$$\mathbb{P}_{p(Z)} \left[ \forall x \in \mathcal{X} \cdot (f_{\hat{\theta}(Z)}(x) - f_{\theta^*}(x))^2 \leq \frac{2K^2\epsilon}{\alpha} \right] \geq 1 - \delta.$$

*Proof.* By Lemma 2, and since  $\nabla_\phi \tilde{L}(g(\theta^*)) = 0$ ,

$$L_p(\hat{\theta}(Z)) - L_p(\theta^*) = \tilde{L}_p(g(\hat{\theta}(Z))) - \tilde{L}_p(g(\theta^*)) \geq \alpha \|g(\hat{\theta}(Z)) - g(\theta^*)\|_F^2. \quad (8)$$

Next, by Lemma 4 and the fact that  $\hat{\theta}$  minimizes  $\hat{L}(\theta; Z)$ ,

$$L_p(\hat{\theta}(Z)) \leq \hat{L}(\hat{\theta}; Z) + \epsilon - \sigma(Z) \leq \hat{L}(\theta^*; Z) + \epsilon - \sigma(Z) \leq L_p(\theta^*) + 2\epsilon \quad (9)$$

with probability at least  $1 - \delta$ . Combining (8) and (9),

$$\|g(\hat{\theta}(Z)) - g(\theta^*)\|_F \leq \sqrt{\frac{2\epsilon}{\alpha}}$$

with probability at least  $1 - \delta$ . Finally, by Lemma 3,

$$(f_{\hat{\theta}(Z)}(x) - f_{\theta^*}(x))^2 = (\tilde{f}_{g(\hat{\theta}(Z))}(x) - \tilde{f}_{g(\theta^*)}(x))^2 \leq K^2 \|g(\hat{\theta}(Z)) - g(\theta^*)\|_2^2 \leq \frac{2K^2\epsilon}{\alpha} \quad (\forall x \in \mathcal{X})$$

with probability at least  $1 - \delta$ , as claimed.  $\square$

As a result, we provide a robust generalization error bound for QNNs with potential distribution shifts.

**Corollary 6.** *Under Assumptions A1 & A2, for any distribution  $q(x)$  with support on  $B_2(0, x_{max})$ ,*

$$\mathbb{P}_{p(Z)} \left[ L_q(\hat{\theta}(Z)) \leq \frac{2K^2\epsilon}{\alpha} \right] \geq 1 - \delta.$$

Finally, we also prove that gradient descent can find the global minima of  $\hat{L}(\theta; Z)$ , which ensures that gradient descent can perform function identification in practice; we give a proof in Appendix C.5.

**Proposition 7.** *All local minima of  $\hat{L}(\theta; Z)$  are global.*

## 4 Quadratic Neural Bandits

A key application of robust generalization bounds is to parametric bandits; this is because, in bandit learning, the distribution of inputs  $x$  used to estimate  $\hat{\theta} \approx \theta^*$  can differ from the distribution under which  $f_{\hat{\theta}}$  is used. Thus, generalization bounds based on notions such as Rademacher complexity cannot be used. Unlike prior literature in bandits, we consider an over-parameterized function that does *not* admit a unique solution; in contrast, recent work on neural tangent kernel bandits [51] assumes that there is a unique, identifiable solution. Note that this assumption cannot hold for quadratic neural networks because they are invariant to transformations such as rotations.

We consider a standard linear bandit [39, 1] with a fixed horizon  $T \in \mathbb{N}$ , but where the expected reward is parameterized by a quadratic neural network instead of a linear function. At each time step  $t$ , the algorithm chooses among a continuum of actions  $x_t \in \mathcal{X}$ , and receives a reward

$$y_t = f_{\theta^*}(x_t) + \xi_t = \sum_{j=1}^k \langle \theta_j^*, x_t \rangle^2 + \xi_t, \quad (10)$$

where  $\theta^* \in \mathbb{R}^{d \times k}$  is an unknown parameter matrix, and  $\xi_t$  are bounded i.i.d. random variables. For simplicity, we assume that  $\mathcal{X} = B_2(0, 1)$  is the unit ball. Then, our goal is to bound the *regret*

$$R(T) = \sum_{t=1}^T (\mathbb{E}_{p(\xi_t)}[y_t] - y^*) \quad \text{where} \quad y^* = \max_{x \in \mathcal{X}} f_{\theta^*}(x).$$

We make the following assumption, which says that  $\phi^* = \theta^* \theta^{*\top}$  has a gap in its top eigenvalue:

**Assumption A3.** Let  $\phi^* = \theta^* \theta^{*\top}$ , and let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$  be the eigenvalues of  $\phi^*$ . There exists a constant  $M \in \mathbb{R}_{>0}$  such that  $\lambda_1 - \lambda_2 \geq 4/M$ .

This assumption ensures that the eigenvectors of  $\phi^*$  to be stable under perturbations. The eigenvectors of  $\phi^*$  correspond to the optimal action  $x^* = \arg \max_{x \in \mathcal{X}} f_{\theta^*}(x)$  since  $f_{\theta^*}(x) = x^\top \phi^* x$ ; thus, it ensures that if  $\hat{\theta} \approx \theta^*$ , then the optimal action  $\hat{x} = \arg \max_{x \in \mathcal{X}} f_{\hat{\theta}}(x)$  satisfies  $\hat{x} \approx x^*$ .

Next, we describe our algorithm, summarized in Algorithm 1. We consider an explore-then-commit strategy for simplicity, since it already achieves the asymptotically optimal regret rate [39]; our approach can similarly be applied to more sophisticated algorithms such as UCB [1] and Thompson sampling [2]. Our algorithm proceeds in two stages: (i) the *exploration stage* (for  $t \in \{1, \dots, m\}$ ), and (ii) the *exploitation stage* (for  $t \in \{m+1, \dots, T\}$ ), where

$$m = \left\lceil \left( \frac{135M(\ell_{\max} + \xi_{\max})^2 d^3 T \sqrt{\log \left( 3 + \frac{\phi_{\max} K T}{\ell_{\max}^2} \right)}}{\phi_{\max}} \right)^{2/3} \right\rceil. \quad (11)$$

In the exploration stage, we randomly choose actions  $x_t \sim p$ , where

$$p(x) = \prod_{i=1}^d \text{Uniform} \left( x_i; \left[ -\frac{1}{\sqrt{d}}, \frac{1}{\sqrt{d}} \right] \right). \quad (12)$$

Note that  $\|x\|_2 \leq 1$  for  $x$  in the support of  $p$ , so  $x_t \in \mathcal{X}$ . Following the discussion in Section 3, for this choice of  $p$ , Assumption A2 holds for the dataset  $Z$  with  $\alpha = 4/(45d^2)$ .

---

**Algorithm 1** Explore-Then-Commit Algorithm for Quadratic Neural Network Bandit

---

**procedure** QUADRATICNEURALBANDITInitialize  $Z \leftarrow \emptyset$ Let  $m$  be as in (11)**for**  $t \in \{1, \dots, m\}$  **do**    Sample i.i.d. action  $x_t \sim p$ , where  $p$  is as in (12)    Take action  $x_t$  and obtain reward  $y_t$  as in (10)    Update  $Z \leftarrow Z \cup \{(x_t, y_t)\}$ **end for**Compute  $\hat{\theta} = \arg \min_{\theta} \hat{L}(\theta; Z)$ , where  $\hat{L}(\theta; Z) = m^{-1} \sum_{i=1}^m (f_{\theta}(x_i) - y_i)^2$ Compute  $\hat{x} = \arg \max_{x \in \mathcal{X}} f_{\hat{\theta}}(x)$ **for**  $t \in \{m+1, \dots, T\}$  **do**    Take action  $x_t = \hat{x}$  and obtain reward  $y_t$  as in (10)**end for****end procedure**

---

Next, we compute an estimate  $\hat{\theta}$  of  $\theta^*$  based on the data  $Z$  collected so far, and compute the optimal action  $\hat{x}$  assuming  $\hat{\theta}$  are the true parameters. Then, in the exploitation stage, we always use action  $\hat{x}$ .

The key challenge providing theoretical guarantees using traditional generalization bounds is handling the optimization problem over  $x \in \mathcal{X}$  used to compute  $\hat{x}$ . Since  $\hat{x}$  is not sampled from the distribution  $p$ , traditional bounds do not provide any guarantees about the accuracy of  $f_{\hat{\theta}}(\hat{x})$  compared to  $f_{\theta^*}(\hat{x})$ . In contrast, Theorem 5 provides a uniform guarantee, so it can be used to bound the regret.

**Theorem 8.** *Under Assumptions A1, A2 & A3, the expected regret of Algorithm 1 is*

$$R(T) \leq C_0 + C_1 \cdot T^{2/3} \left( \log \left( 3 + \frac{8\phi_{\max}KT}{\ell_{\max}^2} \right) \right)^{1/3},$$

where  $C_0$  and  $C_1$  do not depend on  $T$  (see Appendix D).

We give a proof in Appendix D. In particular,  $R(T) = \tilde{O}(T^{2/3})$ . This rate is worse than the usual  $\tilde{O}(\sqrt{T})$  regret since Theorem 5 only admits a  $n^{1/4}$  convergence rate.

## 5 Transfer Learning of QNNs

So far, we have considered shifts in the covariate distribution but not in the label distribution. Now, we consider a transfer learning problem where there is additionally a small shift in the labels. In particular, we assume we have *proxy data*  $Z_p \subseteq \mathcal{X} \times \mathcal{Y}$  from the source domain of the form  $y_{p,i} = f_{\theta_p^*}(x_{p,i}) + \xi_{p,i}$  (for  $i \in [n_p]$ ), where  $\theta_p^* \in \Theta$  are the proxy parameters and  $p(x_p)$  is the source covariate distribution, along with *gold data*  $Z_g \subseteq \mathcal{X} \times \mathcal{Y}$  from the target domain of the form  $y_{g,i} = f_{\theta_g^*}(x_{g,i}) + \xi_{g,i}$  (for  $i \in [n_g]$ ), where  $\theta_g^* \in \Theta$  are the gold parameters and  $q(x_g)$  is the target covariate distribution. We are interested in the setting  $n_p \gg n_g$ , and where  $\|\theta_g^* - \theta_p^*\|_F \leq B$  is small.

We consider a two-stage estimator [9] that first computes an estimate of the proxy parameters  $\hat{\theta}_p = \arg \min_{\theta \in \Theta} \hat{L}(\theta; Z_p)$ , and then computes an estimate of the gold parameters in a way that is constrained towards the proxy parameters. First, note that we have

$$\mathbb{P}_{p(Z)} \left[ L_q(\hat{\theta}_p) \leq \frac{2K^2\epsilon_p}{\alpha} \right] \geq 1 - \frac{\delta}{2},$$

where

$$\epsilon_p = \sqrt{\frac{18\ell_{\max}^2(\ell_{\max}^2 + \xi_{\max}^2)(d^2 \max\left\{1, \log\left(1 + \frac{8\phi_{\max}K n_p}{\ell_{\max}^2}\right)\right\} + \log\frac{4}{\delta})}{n_p}}.$$

where we have highlighted the differences from  $\epsilon$  in (7) in red. Next, we make a technical assumption:

**Assumption A4.** For some  $\sigma_0 \in \mathbb{R}_{>0}$ ,  $\sigma_{\min}(\theta_p^*) \geq \sigma_0$ , where  $\sigma_{\min}(\theta)$  is the  $d$ th singular value of  $\theta$ .

Equivalently, the minimum eigenvalue of  $g(\theta_p^*)$  is positive; intuitively, this assumption ensures a good estimate of  $\theta_p^* \theta_p^{*\top}$  implies a good estimate of  $\theta_p^*$  (up to an orthogonal transformation). Then, letting

$$\hat{B} = B + \frac{1}{\sigma_0} \sqrt{\frac{2\epsilon_p}{\alpha}}$$

be an expanded radius to account for error in our estimate of  $\hat{\theta}_p$ , we use the estimator

$$\hat{\theta}_g = \arg \min_{\theta \in B_2(\hat{\theta}_p, \hat{B})} \hat{L}(\theta; Z_g) \quad \text{where} \quad B_2(\hat{\theta}_p, \hat{B}) = \{\theta \in \Theta \mid \|\theta - \hat{\theta}_p\|_F \leq \hat{B}\}.$$

Note that we have assume  $\hat{B}$  is known; in practice, this constraint can be included as an additive regularization term. Intuitively, this formulation mirrors transfer learning algorithms based on fine-tuning—i.e., initializing the parameters to the proxy data  $\theta_p$  and then taking a small number of steps of stochastic gradient descent (SGD) on the gold data  $\theta_g$ . In particular, SGD can be interpreted as  $L_2$  regularization on the parameters [3], so fine-tuning  $L_2$ -regularizes  $\hat{\theta}_g$  towards  $\theta_p$ .

**Theorem 9.** Under Assumptions A1, A2 & A4, for any  $q(x)$  with support on  $B_2(0, x_{\max})$ ,

$$\mathbb{P}_{p(Z)} \left[ L_q(\hat{\theta}_g) \leq \frac{2K^2\epsilon_g}{\alpha} \right] \geq 1 - \delta,$$

where

$$\epsilon_g = \hat{B} \cdot \sqrt{\frac{18K^2(K^2\hat{B}^2 + \xi_{\max}^2)(d^2 \max\left\{1, \log\left(1 + \frac{8\phi_{\max}K n_g}{\ell_{\max}^2}\right)\right\} + \log\frac{4}{\delta})}{n_g}}.$$

Thus, if  $B$  is small and  $n_p$  is large,  $f_{\hat{\theta}_g}$  is accurate even if  $n_g$  is small; we give a proof in Appendix E.

## 6 Function Identification of ReLU Networks

Next, we consider the identifiability of ReLU networks.

### 6.1 Main Result

We consider a ReLU network  $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ , where  $\mathcal{X} \subseteq \mathbb{R}^d$  and  $\mathcal{Y} \subseteq \mathbb{R}$ , given by

$$f_\theta(x) = \sum_{j=1}^k \sigma(\theta_j^\top x),$$

where  $\theta \in \mathbb{R}^{d \times k}$ , and the ReLU  $\sigma(z) = \mathbb{1}(z \geq 0) \cdot z$  is applied componentwise. Furthermore, we consider the case  $\mathcal{X} = S^{d-1}$ , where  $S^{d-1} \subseteq \mathbb{R}^d$  is the unit sphere in  $d$  dimensions, and consider the input distribution  $p = \text{Uniform}(S^{d-1})$ . We make two assumptions about the true parameters  $\theta^*$ .



**Assumption A5.** We have  $\|\theta_i^*\|_2 = 1$  for all  $i \in [k]$ .

**Assumption A6.** There exists  $\alpha_0 \in \mathbb{R}_{>0}$  such that

$$\|\theta_i^* \pm \theta_{i'}^*\|_2 \geq \alpha_0 \quad (\forall i, i' \in [k]).$$

That is, we assume that the components  $\theta_i^*$  are unit vectors that are separated from one another (as well as from the negatives of the other vectors). This assumption is necessary since if two components are close together, they can approximately be combined into a single one, so closer components are harder to identify.

As before, the estimated parameters  $\hat{\theta}(Z)$  minimize the empirical loss  $\hat{L}(\theta; Z) = n^{-1} \sum_{i=1}^n (f_\theta(x_i) - y_i)^2$ , where  $y_i = f_{\theta^*}(x_i) + \xi_i$ . Since the ground truth parameters are unit vectors, we make the same assumption about the estimated parameters  $\theta$ —i.e.,  $\|\theta_i\|_2 = 1$  for all  $i \in [k]$ .

Now, our main result says that ReLU networks satisfying our assumptions satisfy function identification.

**Lemma 10.** *Under Assumptions A5 & A6, for any  $\eta \in \mathbb{R}_{>0}$  satisfying  $\eta \leq (6126d^2k^2)^{-1}$ , if  $\mathbb{E}_{p(x)} [|f_\theta(x) - f_{\theta^*}(x)|] \leq \eta$ , then*

$$|f_\theta(x) - f_{\theta^*}(x)| \leq 20k^2 \sqrt{d^3 \eta} \quad (\forall x \in \mathcal{X}). \quad (13)$$

Intuitively, this key lemma says that parameters with small errors on distribution  $p$  have small errors everywhere; we provide a proof in Section 6.2. With this result, a similar argument as the proof of Theorem 5 shows the following.

**Theorem 11.** *Under Assumptions A5 & A6, for sufficiently large  $n$ , we have*

$$\mathbb{P}_{p(Z)} \left[ L_q(\hat{\theta}(Z)) \leq 400k^4 d^3 (2\epsilon)^{1/2} \right] \geq 1 - \delta,$$

where

$$\epsilon = \sqrt{\frac{18\ell_{max}^2(\ell_{max}^2 + \xi_{max}^2)(dk \max \left\{ 1, \log \left( 1 + \frac{4kKn}{\ell_{max}^2} \right) \right\} + \log \frac{2}{\delta})}{n}},$$

$\ell_{max} = 2k$  and  $K = 4k$ .

In other words, ReLU networks satisfy function identification. We give a proof in Appendix F.4.

## 6.2 Proof of Lemma 10

A key challenge in proving function identification is that  $f_\theta$  is invariant to permutations of the components of  $\theta_j$ . Thus, we need to establish a mapping from  $\theta_j$  to  $\theta_i^*$ . This mapping is determined by the subset of indices

$$J_i^\alpha = \{j \in [k] \mid \|\theta_j - \theta_i^*\|_2 \leq \alpha\},$$

where  $\alpha \in \mathbb{R}_{>0}$  is a hyperparameter to be chosen later. First, we note that as long as  $\alpha \leq \alpha_0/2$ , then the  $J_i^\alpha$  are disjoint—i.e.,  $J_i^\alpha \cap J_{i'}^\alpha = \emptyset$ ; in particular, if  $j \in J_i^\alpha$ , then for any  $i' \in [k]$ , we have

$$\|\theta_j - \theta_{i'}^*\|_2 \geq \|\theta_i^* - \theta_{i'}^*\|_2 - \|\theta_j - \theta_i^*\|_2 \geq \alpha_0 - \alpha \geq \alpha,$$

so  $j \notin J_{i'}^\alpha$ . As a consequence, if each  $J_i^\alpha$  is a singleton set  $J_i^\alpha = \{j\}$ , then the mapping  $i \mapsto j$  is a bijection. Now, we prove our key lemma, which says that as long as the loss is sufficiently small, then this condition holds.

**Lemma 12.** For any  $\epsilon \in \mathbb{R}_{>0}$ , if

$$\mathbb{E}_{p(x)} [|f_\theta(x) - f_{\theta^*}(x)|] \leq \frac{\epsilon^2 \sqrt{\pi}}{2d} - \frac{\pi k \epsilon^3}{\alpha} - 4k \epsilon^3 \sqrt{\pi d},$$

then  $J_i^\alpha = \{j\}$  is a singleton set for each  $i \in [k]$ .

We give a proof in Section 6.3. In particular, taking

$$\alpha = \frac{\pi k \epsilon^3}{\sqrt{\pi \epsilon^2 / (2d) - \eta} - 4k \epsilon^3 / \sqrt{\pi d}}$$

then  $\mathbb{E}_{p(x)} [|f_\theta(x) - f_{\theta^*}(x)|] = \eta$ ; taking  $\epsilon = \sqrt{6d\eta/\pi^{1/2}}$ , we have

$$\alpha = \frac{\pi k (6d\eta/\pi^{1/2})^{3/2}}{3\eta - \eta - 4k \epsilon^3 / \sqrt{\pi d}} \leq 20k \sqrt{d^3 \eta},$$

where we have used the fact that  $4k \epsilon^3 / \sqrt{\pi d} \leq \eta$  by our assumption on  $\eta$ . Finally, by Lemma 12, each  $J_i^\alpha = \{j\}$  is a singleton set, which implies there exists a bijective map  $i \mapsto j$  such that  $\|\theta_j - \theta_i^*\|_2 \leq \alpha$ ; without loss of generality, assume this map is the identity. Then, we have

$$|f_\theta(x) - f_{\theta^*}(x)| \leq \sum_{i=1}^k |\sigma(\theta_i^\top x) - \sigma(\theta_i^{*\top} x)| \leq k\alpha \leq 20k^2 \sqrt{d^3 \eta}.$$

as claimed.  $\square$

### 6.3 Proof of Lemma 12

It suffices to prove the contrapositive—i.e.,  $J_i^\alpha$  is not a singleton set for some  $i \in [k]$ , then

$$\mathbb{E}_{p(x)} [|f_\theta(x) - f_{\theta^*}(x)|] \geq \frac{\sqrt{\pi \epsilon^2}}{2d} - \frac{\pi k \epsilon^3}{\alpha} - 4k \epsilon^3 \sqrt{\pi d}.$$

In this case, by the pigeonhole principle, there exists  $i \in [k]$  such that  $J_i^\alpha = \emptyset$ . For this  $i$ , define the region

$$X_i = \{x \in \mathcal{X} \mid |\theta_i^{*\top} x| \leq \epsilon\},$$

where  $\epsilon$  is the given hyperparameter. Roughly speaking, our strategy is to show that all the components  $g(\theta_{i'}^{*\top} x)$  (for all  $i' \in [k]$  such that  $i' \neq i$ ) and  $g(\theta_j^\top x)$  (for all  $j \in [k]$ ) are linear in a large fraction of  $X_i$ . Then, since the component  $g(\theta_i^{*\top} x)$  is nonlinear, the gap between it and the remaining components must be large. To be precise, note that

$$f_\theta(x) - f_{\theta^*}(x) = \underbrace{\left( \sum_{j=1}^k g(\theta_j^\top x) - \sum_{i \in [k] \setminus \{i'\}} g(\theta_{i'}^{*\top} x) \right)}_{:=h_i(x)} - g(\theta_i^{*\top} x), \quad (14)$$

where the first term  $h_i(x)$  is linear on a region  $\tilde{X}_i \subseteq X_i$  that is a large fraction of  $X_i$  (in terms of  $p(x)$ ), but the second term  $g(\theta_i^{*\top} x)$  is nonlinear on  $\tilde{X}_i$ . Thus, we can establish a lower bound on the loss  $\mathbb{E}_{p(x)} [|f_\theta(x) - f_{\theta^*}(x)|]$  on  $\tilde{X}_i$ .

One subtlety is that to establish this lower bound,  $\tilde{X}_i$  must be symmetric around the hyperplane  $\theta_i^{*\top} x = 0$ ; thus, our proof approximates the sphere slice  $X_i$  as a cylinder  $X'_i$ , and then cuts out the portions of  $X'_i$  where  $h_i(x)$  may be nonlinear, to obtain  $\tilde{X}'_i$ . Then, we use the above argument with  $\tilde{X}'_i$  instead of  $X'_i$ .

Now, we formalize this argument. First, without loss of generality, we can consider a coordinate system where

$$\theta_i^* = [1 \quad 0 \quad \dots \quad 0]^\top.$$

In this coordinate system, we have

$$X_i = \{x \in \mathcal{X} \mid |x_1| \leq \epsilon\}.$$

Before we continue, we first approximate the set  $X_i$ ; this step is necessary to establish the symmetry property mentioned above. In particular, let

$$\begin{aligned} X'_i &= \{\phi(x) \mid x \in X_i\} \\ \phi(x) &= \left[ x_1 \quad \frac{x_2}{\sqrt{1-x_1^2}} \quad \dots \quad \frac{x_d}{\sqrt{1-x_1^2}} \right]^\top. \end{aligned}$$

Intuitively,  $X'_i$  is a cylinder approximating the slice  $X_i$  of the sphere  $S^{d-1}$  (so its axis is  $\theta_i^*$ )—i.e., it keeps the first component  $x_1$  but projects the remaining ones to form a  $d-2$  sphere. In particular, we can express  $X'_i$  as

$$X'_i = [-\epsilon, \epsilon] \times Z \quad \text{where} \quad Z = S^{d-2},$$

i.e., the product of the interval  $x_1 \in [-\epsilon, \epsilon]$  with the  $d-2$  sphere  $Z = \{z \in \mathbb{R}^{d-1} \mid \|z\|_2 = 1\}$ . This decomposition of  $X'_i$  is important since we will lower bound the loss on each interval  $[-\epsilon, \epsilon] \times \{z\}$  for  $z \in Z$  independently; each of these intervals satisfies the symmetry property described above. Before describing this step, we first show that the loss on  $X'_i$  is a good approximation of the loss on  $X_i$ .

**Lemma 13.** *We have*

$$\left| \int_{X_i} |f_\theta(x) - f_{\theta^*}(x)| dx - \int_{X'_i} |f_\theta(x) - f_{\theta^*}(x)| dx \right| \leq 2k\epsilon^3 \sqrt{d} \cdot |S^{d-2}|.$$

We give a proof in Appendix F.1. By this lemma, it suffices to lower bound the loss on  $X'_i$  instead of  $X_i$ . Now, for each  $z \in Z$ , define the interval

$$X_i^z = [-\epsilon, \epsilon] \times \{z\}.$$

Note that these intervals partition  $X_i$ ; thus, we can lower bound the loss independently on each interval, and then integrate the bound over all intervals  $z \in Z$ . Intuitively, we can lower bound the loss on a given  $X_i^z$  as long as the first term in  $f_\theta(x) - f_{\theta^*}(x)$  shown in (14) is linear on all of  $X_i^z$ .

Thus, the remainder of our proof is divided into two steps: (i) lower bound the loss if this term is linear on  $X_i^z$ , and (ii) upper bound the fraction of  $z \in Z$  for which this term is nonlinear. For the first step, we have the following result.

**Lemma 14.** *For any  $z \in Z$  and  $\beta_0, \beta_1 \in \mathbb{R}$ , we have*

$$\int_{-\epsilon}^{\epsilon} |(\beta_0 + \beta_1 w) - g(w)| dw \geq \frac{\epsilon^2}{4}.$$

We give a proof in Appendix F.2. For the second step, the condition that the first term of (14) is linear on  $X_i^z$  holds if all the terms  $g(\theta_j^\top x)$  (for  $j \in [k]$ ) and  $g(\theta_{i'}^*{}^\top x)$  (for  $i' \in [k] \setminus \{i\}$ ) are linear on  $X_i^z$ .

Thus, it suffices to bound fraction of  $z$  such that each of these terms is nonlinear separately. To this end, note that the ReLU function  $g(w)$  is linear on a region as long as  $w \neq 0$  on that region. Thus, it suffices to *omit* the regions

$$Z_i^\beta = \left\{ z \in Z \mid \exists x_1 \in [-\epsilon, \epsilon] \cdot \beta^\top([x_1] \circ z) = 0 \right\}$$

$$[x_1] \circ z = [x_1 \quad z_1 \quad \dots \quad z_{d-1}]^\top,$$

for all  $\beta = \theta_j$  (for  $j \in [k]$ ) and  $\beta = \theta_{i'}^*$  (for  $i' \in [k] \setminus \{i\}$ ). Then, defining

$$\tilde{Z}_i = Z \setminus \left( \bigcup_{j=1}^k Z_i^{\theta_j} \cup \bigcup_{i' \in [k] \setminus \{i\}} Z_i^{\theta_{i'}^*} \right),$$

we know that the first term  $h_i(x)$  in (14) is linear on  $X_i^z$  for all  $z \in \tilde{Z}_i$ . Our next result bounds the size of  $Z_i^\beta$ .

**Lemma 15.** *For any  $\beta \in \mathbb{R}^d$  such that  $\|\beta\|_2 = 1$  and  $\|\beta \pm \theta_i^*\|_2 \geq \alpha$ , we have*

$$|Z_i^\beta| \leq \frac{2\epsilon \cdot |S^{d-3}|}{\alpha}.$$

We give a proof in Appendix F.3. Furthermore, by Lemma 15 and the fact that  $Z = S^{d-2}$ , we have

$$|\tilde{Z}_i| \geq |S^{d-2}| - \frac{2k\epsilon \cdot |S^{d-3}|}{\alpha}.$$

Now, we can put these results together to prove Lemma 12. First, letting  $\tilde{X}_i' = [-\epsilon, \epsilon] \times \tilde{Z}_i$ , we have

$$\begin{aligned} & \mathbb{E}_{p(x)}[|f_\theta(x) - f_{\theta^*}(x)|] \\ & \geq \frac{1}{|S^{d-1}|} \int_{X_i} |f_\theta(x) - f_{\theta^*}(x)| dx \\ & \geq \frac{1}{|S^{d-1}|} \int_{X_i'} |f_\theta(x) - f_{\theta^*}(x)| dx - \frac{2k\epsilon^3 \sqrt{d} \cdot |S^{d-2}|}{|S^{d-1}|} \\ & \geq \frac{1}{|S^{d-1}|} \int_{\tilde{X}_i'} |f_\theta(x) - f_{\theta^*}(x)| dx - \frac{2k\epsilon^3 \sqrt{d} \cdot |S^{d-2}|}{|S^{d-1}|}, \end{aligned}$$

where the first step follows since  $X_i \subseteq \mathcal{X}$ , the second follows by Lemma 13, and the third follows since  $\tilde{X}_i' \subseteq X_i'$ . Next, letting  $x = [x_1] \circ z$ , we have

$$\begin{aligned} & f_\theta([x_1] \circ z) - f_{\theta^*}([x_1] \circ z) \\ & = h_i([x_1] \circ z) - g(\theta_i^*{}^\top([x_1] \circ z)) \\ & = \beta(z)^\top([x_1] \circ z) - g(\theta_i^*{}^\top([x_1] \circ z)), \end{aligned}$$

where the second equality follows for some  $\beta(z) \in \mathbb{R}^d$  since  $h_i(x)$  is linear on  $\tilde{X}_i^z$  for each  $z \in \tilde{Z}_i$ . Letting  $\beta(z) = [\beta_1(z)] \circ \beta'(z)$ , and letting  $\beta_0(z) = \beta'(z)^\top z$ , we have

$$f_\theta([x_1] \circ z) - f_{\theta^*}([x_1] \circ z) = \beta_0(z) + \beta_1(z)x_1 - g(x_1).$$

Thus, by Lemma 14, we have

$$\begin{aligned}
& \int_{\tilde{X}'_i} |f_\theta(x) - f_{\theta^*}(x)| dx \\
&= \int_{\tilde{Z}_i} \int_{-\epsilon}^{\epsilon} |f_\theta([x_1] \circ z) - f_{\theta^*}([x_1] \circ z)| dx_1 dz \\
&\geq \int_{\tilde{Z}_i} \frac{\epsilon^2}{4} dz \\
&\geq \left( |S^{d-2}| - \frac{2k\epsilon \cdot |S^{d-3}|}{\alpha} \right) \frac{\epsilon^2}{4}.
\end{aligned}$$

As a consequence, we have

$$\begin{aligned}
& \mathbb{E}_{p(x)} [|f_\theta(x) - f_{\theta^*}(x)|] \\
&\geq \left( \frac{|S^{d-2}|}{|S^{d-1}|} - \frac{2k\epsilon}{\alpha} \cdot \frac{|S^{d-3}|}{|S^{d-1}|} \right) \frac{\epsilon^2}{4} - 2k\epsilon^3 \sqrt{d} \cdot \frac{|S^{d-2}|}{|S^{d-1}|}.
\end{aligned}$$

Finally, for any  $d' \leq d$ , note that

$$\frac{2\sqrt{\pi}}{d} \leq \frac{|S^{d'-1}|}{|S^{d'-2}|} \leq 2\sqrt{\pi},$$

which follows from the formula of volume of the  $n$ -sphere  $|S^n| = \pi^{n/2}/\Gamma((n/2) + 1)$  and the Gamma function identity  $\Gamma(w + 1) = w\Gamma(w)$ .  $\square$

## 7 Conclusion

We have presented results demonstrating that over-parameterization does not fundamentally harm learning models that are robust to arbitrary distribution shifts. In particular, even though we can no longer identify the true parameters for QNNs, we show that we can identify the true function, thereby enabling us to prove new results in bandits and transfer learning. Finally, we also prove function identification for a subclass of ReLU networks.

## References

- [1] Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári. Improved algorithms for linear stochastic bandits. In *NIPS*, volume 11, pages 2312–2320, 2011.
- [2] S. Agrawal and N. Goyal. Thompson sampling for contextual bandits with linear payoffs. In *International Conference on Machine Learning*, pages 127–135. PMLR, 2013.
- [3] A. Ali, E. Dobriban, and R. Tibshirani. The implicit regularization of stochastic gradient flow for least squares. In *International Conference on Machine Learning*, pages 233–244. PMLR, 2020.
- [4] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein. Neural module networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 39–48, 2016.
- [5] R. Arora, A. Basu, P. Mianjy, and A. Mukherjee. Understanding deep neural networks with rectified linear units. In *ICLR*, 2018.

- [6] S. Arora, R. Ge, B. Neyshabur, and Y. Zhang. Stronger generalization bounds for deep nets via a compression approach. In *International Conference on Machine Learning*, pages 254–263. PMLR, 2018.
- [7] F. Bach, J. Mairal, and J. Ponce. Convex sparse matrix factorizations. *arXiv preprint arXiv:0812.1869*, 2008.
- [8] P. L. Bartlett and S. Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [9] H. Bastani. Predicting with proxies: Transfer learning in high dimension. *Management Science*, 2020.
- [10] S. Ben-David, J. Blitzer, K. Crammer, and F. Pereira. Analysis of representations for domain adaptation. In *Advances in neural information processing systems*, pages 137–144, 2007.
- [11] J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. Wortman. Learning bounds for domain adaptation. In *Advances in neural information processing systems*, pages 129–136, 2008.
- [12] E. J. Candes and Y. Plan. Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *IEEE Transactions on Information Theory*, 57(4): 2342–2359, 2011.
- [13] J. Cohen, E. Rosenfeld, and Z. Kolter. Certified adversarial robustness via randomized smoothing. In *International Conference on Machine Learning*, pages 1310–1320. PMLR, 2019.
- [14] I. Diakonikolas, G. Kamath, D. Kane, J. Li, A. Moitra, and A. Stewart. Robust estimators in high-dimensions without the computational intractability. *SIAM Journal on Computing*, 48(2): 742–864, 2019.
- [15] S. Du and J. Lee. On the power of over-parametrization in neural networks with quadratic activation. In *International Conference on Machine Learning*, pages 1329–1338. PMLR, 2018.
- [16] S. Du, J. Lee, H. Li, L. Wang, and X. Zhai. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pages 1675–1685. PMLR, 2019.
- [17] J. Duchi and H. Namkoong. Learning models with uniform performance via distributionally robust optimization. *arXiv preprint arXiv:1810.08750*, 2018.
- [18] D. Foster and A. Rakhlin. Beyond ucb: Optimal and efficient contextual bandits with regression oracles. In *International Conference on Machine Learning*, pages 3199–3210. PMLR, 2020.
- [19] R. Gao, T. Cai, H. Li, C.-J. Hsieh, L. Wang, and J. D. Lee. Convergence of adversarial training in overparametrized neural networks. *Advances in Neural Information Processing Systems*, 32: 13029–13040, 2019.
- [20] R. Ge, C. Jin, and Y. Zheng. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. In *International Conference on Machine Learning*, pages 1233–1242. PMLR, 2017.
- [21] R. Ge, J. D. Lee, and T. Ma. Learning one-hidden-layer neural networks with landscape design. *arXiv preprint arXiv:1711.00501*, 2017.

- [22] S. Goel and A. R. Klivans. Learning neural networks with two nonlinear layers in polynomial time. In *COLT*, pages 1470–1499, 2019.
- [23] S. Goel, A. Klivans, P. Manurangsi, and D. Reichman. Tight hardness results for training depth-2 relu networks. In *ICTS*, 2021.
- [24] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [25] D. Hendrycks and T. Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- [26] D. Hsu, S. M. Kakade, and P. Liang. Identifiability and unmixing of latent parse trees. *arXiv preprint arXiv:1206.3137*, 2012.
- [27] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry. Adversarial examples are not bugs, they are features. *arXiv preprint arXiv:1905.02175*, 2019.
- [28] A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *arXiv preprint arXiv:1806.07572*, 2018.
- [29] M. Kabanava, R. Kueng, H. Rauhut, and U. Terstiege. Stable low-rank matrix recovery via null space properties. *Information and Inference: A Journal of the IMA*, 5(4):405–441, 2016.
- [30] M. J. Kearns, U. V. Vazirani, and U. Vazirani. *An introduction to computational learning theory*. MIT press, 1994.
- [31] P. W. Koh, S. Sagawa, H. Marklund, S. M. Xie, M. Zhang, A. Balsubramani, W. Hu, M. Yasunaga, R. L. Phillips, I. Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. *arXiv preprint arXiv:2012.07421*, 2020.
- [32] B. Lake and M. Baroni. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks. In *International conference on machine learning*, pages 2873–2882. PMLR, 2018.
- [33] Y. Li, T. Ma, and H. Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In *Conference On Learning Theory*, pages 2–47. PMLR, 2018.
- [34] P. M. Long and H. Sedghi. Generalization bounds for deep convolutional neural networks. *arXiv preprint arXiv:1905.12600*, 2019.
- [35] B. Neyshabur, S. Bhojanapalli, D. McAllester, and N. Srebro. Exploring generalization in deep learning. *arXiv preprint arXiv:1706.08947*, 2017.
- [36] A. Raghunathan, J. Steinhardt, and P. Liang. Certified defenses against adversarial examples. *arXiv preprint arXiv:1801.09344*, 2018.
- [37] M. T. Ribeiro, T. Wu, C. Guestrin, and S. Singh. Beyond accuracy: Behavioral testing of nlp models with checklist. *arXiv preprint arXiv:2005.04118*, 2020.
- [38] L. Ruis, J. Andreas, M. Baroni, D. Bouchacourt, and B. M. Lake. A benchmark for systematic generalization in grounded language understanding. *arXiv preprint arXiv:2003.05161*, 2020.

- [39] P. Rusmevichientong and J. N. Tsitsiklis. Linearly parameterized bandits. *Mathematics of Operations Research*, 35(2):395–411, 2010.
- [40] H. Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of statistical planning and inference*, 90(2):227–244, 2000.
- [41] M. Soltanolkotabi, A. Javanmard, and J. D. Lee. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information Theory*, 65(2):742–769, 2018.
- [42] J. Steinhardt, P. W. Koh, and P. Liang. Certified defenses for data poisoning attacks. *arXiv preprint arXiv:1706.03691*, 2017.
- [43] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [44] R. Taori, A. Dave, V. Shankar, N. Carlini, B. Recht, and L. Schmidt. Measuring robustness to natural distribution shifts in image classification. *arXiv preprint arXiv:2007.00644*, 2020.
- [45] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.
- [46] R. Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [47] M. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*. Book Draft (Working Publication), 2016. URL [https://www.stat.berkeley.edu/~wainwrig/nachdiplom/Chap2\\_Sep10\\_2015.pdf](https://www.stat.berkeley.edu/~wainwrig/nachdiplom/Chap2_Sep10_2015.pdf).
- [48] J. Wen, C.-N. Yu, and R. Greiner. Robust learning under uncertain test distributions: Relating covariate shift to model misspecification. In *International Conference on Machine Learning*, pages 631–639. PMLR, 2014.
- [49] K. Xu, X. Zhao, H. Bastani, and O. Bastani. Group-sparse matrix factorization for transfer learning of word embeddings. In *International Conference on Machine Learning*, 2021.
- [50] Y. Yu, T. Wang, and R. J. Samworth. A useful variant of the davis–kahan theorem for statisticians. *Biometrika*, 102(2):315–323, 2015.
- [51] D. Zhou, L. Li, and Q. Gu. Neural contextual bandits with ucb-based exploration. In *International Conference on Machine Learning*, pages 11492–11502. PMLR, 2020.



## A Additional Related Work

**Low-rank matrix factorization.** Our notion of functional identification for quadratic neural network is related to the low-rank matrix factorization literature. However, they impose a low-rank structure on the matrix to recover, and hence typically require extra conditions to identify the matrix—e.g., the restricted isometry property (RIP) [12, 20], or bounded  $\ell_2$  norm of noise vector [29]. In contrast, we consider a more general case and do not assume any underlying structure of the matrix; in particular, since our goal is to capture over-parameterization of neural networks, our matrix is usually decomposed as  $\phi = \theta\theta^\top$ , where  $\theta \in \mathbb{R}^{d \times k}$  and  $k \geq d$  (and  $\phi$  is not necessarily low-rank). Also, we study the prediction error of neural networks in the presence of distribution shifts, whereas the goal of the low-rank literature is to recover the true matrix.

**Multi-armed bandits.** Prior literature on parameterized bandits has considered a number of functional forms, ranging from linear [1, 39] to neural tangent kernels [51]. Most of this work makes a *realizability* assumption that the model family contains the true model;<sup>2</sup> implicitly, they consider model families where there is a *unique*, identifiable true parameter. These assumptions are necessary precisely due to the fact that the test and training distributions are different; thus, much of the bandit literature has focused on proving parameter identification results to enable learning. In contrast, the identifiability assumption does not hold for quadratic neural networks because they are invariant to parameter transformations. To the best of our knowledge, we consider the first over-parameterized bandit problem that considers a model that is not parameter-identifiable; we find that similar regret results hold as long as the function represented by the model can be identified. Separately, [18] makes a general connection between online regression oracles and the regret of a bandit algorithm; however, their approach only provides good guarantees when the regression oracle returns a model that generalizes off-distribution. Finally, recent work on UCB with neural tangent kernels [51] provides general regret bounds, but their bound is only sublinear under conditions such as the true reward function having small RKHS norm (see Remark 4.8 in their paper), which amounts to assuming they can recover the true parameters.

**Novelty.** We briefly discuss the novelty of our results compared to existing work. First, to the best of our knowledge, all our results for ReLU networks in Section 6 are novel. For QNNs, the results in Section 3 are novel. To the best of our knowledge, the proof strategy in our main result, Theorem 5, is novel, though we note that the preceding lemmas are based on standard arguments—e.g., bounding the convexity of  $\tilde{L}_p(\phi)$  (Lemma 2) and the Lipschitz constant (Lemma 3) of  $\tilde{f}_\phi$ ; also, Lemma 4 relies on a standard covering number argument. For our applications to bandits and transfer learning, our key novel results are Lemma 26 for bandits, which proves smoothed bounded response for quadratic neural networks, and Lemma 27 for transfer learning. Finally, to the best of our knowledge, our arguments in Section B are novel.

## B Generalization Bounds for Neural Module Networks

While function identification enables robust generalization, many data generating processes are too complex to be identifiable. Neural module networks are designed to break complex prediction problems into smaller tasks that are individually easier to solve. These models take two kinds of input: (i) a sequence of tokens  $w$  (e.g., word embeddings) indicating the correct composition of modules, and (ii) the input  $x$  to the modules. Then, the model predicts the sequence of modules  $j_1 \dots j_T$  based on  $w$ , and runs the modules in sequence to obtain output  $x' = f_{j_T}(\dots(f_{j_1}(x))\dots)$ .

---

<sup>2</sup>Slightly different from realizability in PAC learning [30], which says there is a model with zero true loss.

We study conditions under which neural module networks can robustly generalize. Rather than study arbitrary distribution shifts, we consider two separate shifts:

- **Module inputs:** We assume that the individual modules are identifiable; as a consequence, we assume the shift to the module input  $x$  can be arbitrary.
- **Module composition:** We consider shifts to the token sequence  $w$ . If the model mapping  $w$  to  $j_1 \dots j_T$  is identifiable, then the entire model is identifiable. Instead, we show that when this model is not identifiable, compositional structure can still aid generalization. Intuitively, we show that while small shifts in the compositional structure can cause large shifts in the distribution  $p(w)$ , models that leverage the structure of  $p(w)$  can still generalize well.

In more detail, consider a simplified neural module network  $f$ , which includes (i) a set of neural modules  $\{f_j : \mathcal{X} \rightarrow \mathcal{X}\}_{j=1}^k$ , and (ii) a parser  $g : \mathcal{Z}^T \rightarrow [k]^T$ , where  $\mathcal{Z} \subseteq \mathbb{R}^r$ , with model class  $g \in \mathcal{G}$ . We assume each component of  $f_j(x)$  is computed by a separate quadratic neural network; we discuss the architecture of  $g$  below. Then, given an input  $x \in \mathcal{X} \subseteq \mathbb{R}^d$  and  $w \in \mathcal{W} = \mathcal{Z}^T$ , the corresponding neural module network  $f : \mathcal{X} \times \mathcal{W} \rightarrow \mathcal{X}$  is defined by

$$f(x, w) = (f_{j_T} \circ \dots \circ f_{j_1})(x) = f_{j_T}(\dots(f_{j_1}(x))\dots) \quad \text{where} \quad j_1 \dots j_T = g(w).$$

We assume that  $g$  has compositional structure—i.e., for some  $\tilde{g} : [k] \times \mathcal{Z} \rightarrow [k]$ , we have

$$g(w) = j_1 \dots j_T \quad \text{where} \quad j_t = \begin{cases} 0 & \text{if } t = 0 \\ \tilde{g}(z_t, j_{t-1}) & \text{otherwise,} \end{cases}$$

where  $w = z_1 \dots z_T$ . Intuitively,  $w$  is a sequence of word vectors; then, the current neural module  $j_t = \tilde{g}(z_t, j_{t-1})$  depends both on the current word vector  $z_t$  and the previous neural module  $j_{t-1}$ . First, we assume that the individual modules have been functionally identified.

We assume we have fully labeled data we can use to train the neural modules—i.e., for each input  $x$  and sequence  $w$ , we have both the desired sequence  $j_1 \dots j_T$  of neural modules, as well as the entire execution  $x_0, x_1, \dots, x_T$ , where  $x_0 = x$  and  $x_{t+1} = f_{j_t}(x_t)$  otherwise. Thus, we can use supervised learning to train the neural modules;<sup>3</sup> in particular, we can construct labeled examples  $(j_{t-1}, z_t, j_t)$  used to train the parser  $\tilde{g}$ , and labeled examples  $(x_t, x_{t+1})$  to train the modules  $f_{j_t}$ . For simplicity, we assume we have a uniform lower bound  $n$  on the number of training examples for the parser and for each module. Then, we have the following straightforward result:

**Lemma 16.** *Under Assumptions A1 & A2, with probability at least  $1 - dk\delta$ , for each  $j \in [k]$ ,*

$$\|\hat{f}_j(x) - f_j^*(x)\|_2 \leq \sqrt{\frac{2dK^2\epsilon}{\alpha}} =: \epsilon_f \quad (\forall x \in \mathcal{X}),$$

where  $\hat{f}_j$  is the estimated module and  $f_j^*$  is the ground truth module.

This result follows straightforwardly from Theorem 5 along with a union bound. In contrast, we do not assume the parsing model robustly generalizes, but only on distribution. For the subsequent analysis, we can use any neural network models that satisfy the statement of Lemma 16.

**Lemma 17.** *Under Assumptions A1 & A2, with probability at least  $1 - \delta$ , we have*

$$\mathbb{P}_{\tilde{p}(z,j)} \left[ \hat{g}(z, j) \neq \tilde{g}^*(z, j) \right] \leq 4R_n(\mathcal{G}) + \sqrt{\frac{2 \log(2/\delta)}{n}} =: \epsilon_g,$$

<sup>3</sup>Neural modules are often trained with only partial supervision [4]; we leave an analysis of this strategy to future work since our focus is on understanding generalization rather than learning dynamics.

where  $R_n(\mathcal{G})$  is the Rademacher complexity of  $\mathcal{G}$  (including its loss function), where  $p(z, j) = T^{-1} \sum_{t=1}^T p_t(z, j)$ , and where

$$\tilde{p}_t(z, j) = \begin{cases} \mathbb{1}(j = 0) \cdot \tilde{p}(z) & \text{if } t = 1 \\ \sum_{j'=1}^k \int \mathbb{1}(j = \tilde{g}^*(z', j')) \cdot \tilde{p}(z | z') \cdot \tilde{p}_{t-1}(z', j') dz' & \text{otherwise.} \end{cases}$$

This result is a standard Rademacher generalization bound [8]. Note that we have also assumed that the distribution over token sequences is structured, which is necessary for our compositional implementation of  $g$  to generalize, even on distribution. Intuitively, the distribution over  $(z, j)$  consists of both a unigram model over the word vectors:

$$p(z_1, \dots, z_T) = \prod_{t=1}^T \tilde{p}(z_t | z_{t-1}),$$

where we define  $\tilde{p}(z_1 | z_0) = \tilde{p}(z_1)$ , as well as a unigram model over neural modules:

$$p(j_1 \dots j_T | z_1, \dots, z_T) = \prod_{t=1}^T \mathbb{1}(j_t = \tilde{g}^*(z_t, j_{t-1})).$$

Next, we consider a shifted distribution  $\tilde{q}(z | z')$ , which is close to  $\tilde{p}(z | z')$ .

**Assumption A7.** We have  $\|\tilde{q}(\cdot | z') - \tilde{p}(\cdot | z')\|_{TV} \leq \alpha$ .

Importantly, despite this assumption, the shift between the overall distributions  $p(z_1, \dots, z_T)$  and  $q(z_1, \dots, z_T)$  can still be large since it compounds exponentially across the steps  $t \in [T]$ .

**Proposition 18.** *There exist  $p$  and  $q$  that satisfy Assumption A7, but  $\|p - q\|_{TV} = 2(1 - (1 - \alpha/2)^T)$ .*

That is, even if the single step probabilities  $\tilde{p}(z | z')$  and  $\tilde{q}(z | z')$  have total variation (TV) distance bounded as in Assumption A7, the overall distributions  $p$  and  $q$  can have TV distance exponentially close to the maximum possible distance of 2 in  $T$ ; we give a proof in Appendix G.1.

We show that neural module networks generalize since  $\hat{g}$  leverages the compositional structure of  $p$ . First, we show that under Assumption A7, the overall shift in the input distribution of  $\hat{g}$  is bounded:

**Lemma 19.** *Under Assumptions A1, A2 & A7, we have  $\|\tilde{q} - \tilde{p}\|_{TV} \leq T\alpha$ , where  $\tilde{p}$  is defined in Lemma 17 and  $\tilde{q}$  is defined in Assumption A7.*

That is, while the shift can compound across steps  $t$ , it does so only linearly; we give a proof in Appendix G.2. Next, we show that as a consequence, the error of  $\hat{g}$  is bounded.

**Lemma 20.** *Under Assumptions A1, A2 & A7, and assuming that  $\mathbb{P}_{p(z,j)}[\hat{g}(z, j) \neq \tilde{g}^*(z, j)] \leq \epsilon_g$ , we have that  $\mathbb{P}_{p(w)}[\hat{g}(w) \neq g^*(w)] \leq T\epsilon_g$ .*

We give a proof in Appendix G.3. Finally, we have our main result.

**Theorem 21.** *Under Assumptions A1, A2 & A7, with probability at least  $1 - (dk + 1)\delta$ , we have*

$$\mathbb{P}_{q(w)} \left[ \|\hat{f}(x, w) - f^*(x, w)\|_2 \leq T\epsilon_f \cdot \max\{K^{T-1}, 1\} \right] \geq 1 - T\epsilon_g - T^2\alpha.$$

We give a proof in Appendix G.4. Intuitively, Theorem 21 says that the error of the neural module network is linear in  $T$  as long as  $K \leq 1$ . Note that even if there is no distribution shift, its error is

$$\mathbb{P}_{p(w)} \left[ \|\hat{f}(x, w) - f^*(x, w)\|_2 \leq T\epsilon_f \cdot \max\{K^{T-1}, 1\} \right] \geq 1 - T\epsilon_g,$$

by the same argument as the proof of Theorem 21. The exponential dependence on  $K$  is unavoidable since  $K > 1$  says that the modules  $f_j$  can expand the input, which leads to exponential blowup in the magnitude of the output as a function of  $T$ , which also makes the estimation error exponential in  $T$ . Thus, the only cost to the distribution shift from  $p$  to  $q$  is the additional error probability  $T^2\alpha$ .

## C Proofs for Section 3

### C.1 Proof of Minimum Eigenvalue for Uniform Distribution

In this section, we prove the claim that Assumption A2 holds for the covariate distribution where  $x_i$  is an i.i.d. random variable with distribution  $\text{Uniform}(x_i; [-1/2, 1/2])$ . To this end, note that

$$\begin{aligned} \mathbb{E}_{p(x)}[(x^\top \Delta x)^2] &= \mathbb{E}_{p(x)} \left[ \left( \sum_{i,j=1}^d x_i x_j \Delta_{ij} \right)^2 \right] \\ &= \mathbb{E}_{p(x)} \left[ \sum_i x_i^4 \Delta_{ii}^2 + \sum_{i \neq j} x_i^2 x_j^2 \Delta_{ii} \Delta_{jj} + 2 \sum_{i \neq j} x_i^2 x_j^2 \Delta_{ij}^2 \right] \\ &= \frac{1}{80} \sum_i \Delta_{ii}^2 + \frac{1}{144} \sum_{i \neq j} \Delta_{ii} \Delta_{jj} + \frac{1}{72} \sum_{i \neq j} \Delta_{ij}^2 \\ &= \left( \frac{1}{80} - \frac{1}{144} \right) \sum_i \Delta_{ii}^2 + \frac{1}{144} \left( \sum_i \Delta_{ii} \right)^2 + \frac{1}{72} \sum_{i \neq j} \Delta_{ij}^2 \\ &\geq \frac{1}{180} \|\Delta\|_F^2, \end{aligned}$$

as claimed.  $\square$

### C.2 Proof of Lemma 2

We use the notation  $U : \nabla^2 f(\phi) : V$  to denote the matrix inner product  $\langle U, \nabla^2 f(\phi)(V) \rangle$  for  $U, V \in \mathbb{R}^{d \times d}$ . The Hessian  $\nabla^2 f(\phi)$  can be viewed as a  $d^2 \times d^2$  matrix. As everything here is bounded, we can exchange the expectation and differentiation. Therefore, the Hessian of our loss function has for any symmetric matrix  $\Delta$

$$\Delta : \nabla^2 \tilde{L}_p(\phi) : \Delta = 2\mathbb{E}_{p(x)}[(x^\top \Delta x)^2] \geq 2\alpha \|\Delta\|_F^2,$$

where the last inequality uses Assumption A2.  $\square$

### C.3 Proof of Lemma 3

By our definition, for any  $\phi, \phi' \in \Phi$ ,

$$|\tilde{f}_\phi(x) - \tilde{f}_{\phi'}(x)| = |(x^\top (\phi - \phi') x)| \leq x_{\max}^2 \|\phi - \phi'\|_F.$$

Given our quadratic loss function, we have

$$\begin{aligned}
& |(\tilde{f}_\phi(x) - \tilde{f}_{\phi^*}(x))^2 - (\tilde{f}_{\phi'}(x) - \tilde{f}_{\phi^*}(x))^2| \\
& \leq |\tilde{f}_\phi(x) - \tilde{f}_{\phi^*}(x) + \tilde{f}_{\phi'}(x) - \tilde{f}_{\phi^*}(x)| |\tilde{f}_\phi(x) - \tilde{f}_{\phi'}(x)| \\
& \leq 4\phi_{\max} x_{\max}^4 \|\phi - \phi'\|_F.
\end{aligned}$$

Next, the true loss satisfies

$$|\tilde{L}_p(\phi) - \tilde{L}_p(\phi')| \leq \mathbb{E}_{p(x)} [ |(\tilde{f}_\phi(x) - \tilde{f}_{\phi^*}(x))^2 - (\tilde{f}_{\phi'}(x) - \tilde{f}_{\phi^*}(x))^2| ] \leq 4\phi_{\max} x_{\max}^4 \|\phi - \phi'\|_F.$$

Finally, the empirical loss satisfies

$$\begin{aligned}
|\hat{L}(\phi; Z) - \hat{L}(\phi'; Z)| &= \left| \frac{1}{n} \sum_{i=1}^n [(\tilde{f}_\phi(x_i) - y_i)^2 - (\tilde{f}_{\phi'}(x_i) - y_i)^2] \right| \\
&\leq \frac{1}{n} \sum_{i=1}^n |(\tilde{f}_\phi(x_i) - \tilde{f}_{\phi^*}(x_i))^2 - (\tilde{f}_{\phi'}(x_i) - \tilde{f}_{\phi^*}(x_i))^2| \\
&\quad + \frac{1}{n} \sum_{i=1}^n |\xi_i| \cdot |\tilde{f}_\phi(x_i) - \tilde{f}_{\phi'}(x_i)| \\
&\leq (4\phi_{\max} x_{\max}^4 + 2\xi_{\max} x_{\max}^2) \|\phi - \phi'\|_F,
\end{aligned}$$

as claimed.  $\square$

#### C.4 Proof of Lemma 4

First, we have the following results:

**Lemma 22** (Covering Number of Ball). *For a ball in  $\mathbb{R}^{n_1 \times n_2}$  with radius  $R$  with respect to any norm, there exists an  $\epsilon$ -net  $\mathcal{E}$  such that*

$$|\mathcal{E}| \leq \left(1 + \frac{2R}{\epsilon}\right)^{n_1 n_2}.$$

*Proof.* This claim follows by a direct application of Proposition 4.2.12 in [46].  $\square$

**Lemma 23** (Hoeffding's Inequality for Subgaussian Random Variables). *Letting  $\{z_i\}_{i=1}^n$  be a set of independent  $\sigma$ -subgaussian random variables, then for all  $t \geq 0$ , we have*

$$\Pr \left[ \frac{1}{n} \sum_{i=1}^n z_i \geq t \right] \leq \exp \left( -\frac{2nt^2}{\sigma^2} \right).$$

*Proof.* See Proposition 2.1 of [47].  $\square$

Now, we prove Lemma 4. Consider an  $\epsilon/(4K)$ -net  $\mathcal{E}$  with respect to Frobenius norm. Then, for any  $\phi \in \Phi$ , there exists  $\phi' \in \mathcal{E}$  such that

$$|(\hat{L}(\phi; Z) - \tilde{L}_p(\phi)) - (\hat{L}(\phi'; Z) - \tilde{L}_p(\phi'))| \leq 2K \|\phi - \phi'\|_F \leq \frac{\epsilon}{2}.$$

Therefore, we have

$$\begin{aligned}
& \mathbb{P}_{p(Z)} \left[ \sup_{\theta} |\hat{L}(g(\theta); Z) - L_p(g(\theta)) - \sigma(Z)| \geq \epsilon \right] \\
&= \mathbb{P}_{p(Z)} \left[ \sup_{\phi \in \Phi} |\hat{\tilde{L}}(\phi; Z) - \tilde{L}_p(\phi) - \sigma(Z)| \geq \epsilon \right] \\
&\leq \mathbb{P}_{p(Z)} \left[ \max_{\phi \in \mathcal{E}} |\hat{\tilde{L}}(\phi; Z) - \tilde{L}_p(\phi) - \sigma(Z)| \geq \frac{\epsilon}{2} \right] \\
&\leq \sum_{\phi \in \mathcal{E}} \mathbb{P}_{p(Z)} \left[ |\hat{\tilde{L}}(\phi; Z) - \tilde{L}_p(\phi) - \sigma(Z)| \geq \frac{\epsilon}{2} \right]. \tag{15}
\end{aligned}$$

Now, defining

$$\bar{\tilde{L}}(\phi; Z) = \frac{1}{n} \sum_{i=1}^n (\tilde{f}_{\phi}(x_i) - \tilde{f}_{\phi^*}(x_i))^2 \quad \text{and} \quad \tilde{\eta}(\phi; Z) = \frac{1}{n} \sum_{i=1}^n (\tilde{f}_{\phi}(x_i) - \tilde{f}_{\phi^*}(x_i))\xi_i,$$

and recalling that  $\sigma(Z) = n^{-1} \sum_{i=1}^n \xi_i^2$ , then we have

$$\hat{\tilde{L}}(\phi; Z) = \bar{\tilde{L}}(\phi; Z) + 2\tilde{\eta}(\phi; Z) + \sigma(Z).$$

Thus, continuing from (15), we have

$$\begin{aligned}
& \sum_{\phi \in \mathcal{E}} \mathbb{P}_{p(Z)} \left[ |\hat{\tilde{L}}(\phi; Z) - \tilde{L}_p(\phi) - \sigma(Z)| \geq \frac{\epsilon}{2} \right] \\
&\leq \sum_{\phi \in \mathcal{E}} \mathbb{P}_{p(Z)} \left[ |\bar{\tilde{L}}(\phi; Z) - \tilde{L}_p(\phi)| + 2|\tilde{\eta}(\phi; Z)| \geq \frac{\epsilon}{2} \right] \\
&\leq \sum_{\phi \in \mathcal{E}} \left( \mathbb{P}_{p(Z)} \left[ |\bar{\tilde{L}}(\phi; Z) - \tilde{L}_p(\phi)| \geq \frac{\epsilon}{6} \right] + \mathbb{P}_{p(Z)} \left[ |\tilde{\eta}(\phi; Z)| \geq \frac{\epsilon}{6} \right] \right). \tag{16}
\end{aligned}$$

For the first term in (16), note that  $|\tilde{f}_{\phi}(x) - \tilde{f}_{\phi^*}(x)|^2 \leq \ell_{\max}^2$ , so  $(\tilde{f}_{\phi}(x) - \tilde{f}_{\phi^*}(x))^2$  is  $\ell_{\max}^2$ -subgaussian; thus, by Lemma 23, we have

$$\sum_{\phi \in \mathcal{E}} \mathbb{P}_{p(Z)} \left[ |\bar{\tilde{L}}(\phi; Z) - \tilde{L}_p(\phi)| \geq \frac{\epsilon}{6} \right] \leq 2|\mathcal{E}| \cdot \exp\left(-\frac{n\epsilon^2}{18\ell_{\max}^4}\right). \tag{17}$$

Next, for the second term in (16), note that  $|(\tilde{f}_{\phi}(x_i) - \tilde{f}_{\phi^*}(x_i))\xi_i| \leq \ell_{\max}\xi_{\max}$ , so  $(\tilde{f}_{\phi}(x_i) - \tilde{f}_{\phi^*}(x_i))\xi_i$  is  $\ell_{\max}\xi_{\max}$ -subgaussian; thus, by Lemma 23, we have

$$\sum_{\phi \in \mathcal{E}} \mathbb{P}_{p(Z)} \left[ |\tilde{\eta}(\phi; Z)| \geq \frac{\epsilon}{6} \right] \leq 2|\mathcal{E}| \cdot \exp\left(-\frac{n\epsilon^2}{18\ell_{\max}^2\xi_{\max}^2}\right). \tag{18}$$

Combining (17) & (18), continuing from (16), we have

$$\begin{aligned}
& \sum_{\phi \in \mathcal{E}} \left( \mathbb{P}_{p(Z)} \left[ |\bar{\tilde{L}}(\phi; Z) - \tilde{L}_p(\phi)| \geq \frac{\epsilon}{6} \right] + \mathbb{P}_{p(Z)} \left[ |\tilde{\eta}(\phi; Z)| \geq \frac{\epsilon}{6} \right] \right) \\
&\leq 4|\mathcal{E}| \cdot \exp\left(-\frac{n\epsilon^2}{18\ell_{\max}^2(\ell_{\max}^2 + \xi_{\max}^2)}\right) \\
&\leq 2 \left( 1 + \frac{8\phi_{\max}K}{\epsilon} \right)^{d^2} \cdot \exp\left(-\frac{n\epsilon^2}{18\ell_{\max}^2(\ell_{\max}^2 + \xi_{\max}^2)}\right) \\
&= 2 \exp\left(-\frac{n\epsilon^2}{18\ell_{\max}^2(\ell_{\max}^2 + \xi_{\max}^2)} + d^2 \log\left(1 + \frac{8\phi_{\max}K}{\epsilon}\right)\right), \tag{19}
\end{aligned}$$

where for the first inequality, we have used  $\max\{\ell_{\max}^2, \xi_{\max}^2\} \leq \ell_{\max}^2 + \xi_{\max}^2$ , and the second inequality follows since by Lemma 22, the covering number of the  $\epsilon$ -net  $\mathcal{E}$  of  $\Phi$  satisfies

$$|\mathcal{E}| \leq \left(1 + \frac{2\phi_{\max}}{\epsilon}\right)^{d^2}.$$

Finally, we choose  $\epsilon$  so that (38) is smaller than  $\delta$ —in particular, letting

$$\epsilon = \sqrt{\frac{18\ell_{\max}^2(\ell_{\max}^2 + \xi_{\max}^2)}{n} \left(d^2 \max\left\{1, \log\left(1 + \frac{8\phi_{\max}Kn}{\ell_{\max}^2}\right)\right\} + \log\frac{2}{\delta}\right)}.$$

then continuing (38), we have

$$2 \exp\left(-\frac{n\epsilon^2}{18\ell_{\max}^2(\ell_{\max}^2 + \xi_{\max}^2)} + d^2 \log\left(1 + \frac{8\phi_{\max}K}{\epsilon}\right)\right) \leq \delta,$$

as claimed.  $\square$

## C.5 Proof of Proposition 7

$\hat{L}(\phi; Z)$  is twice differentiable and convex in  $\phi$ . Note that the minimization problem of  $\hat{L}(\theta; Z)$  is equivalent to that of  $\hat{L}(g(\hat{\theta}); Z)$ . We consider two cases. First, consider the case where  $\hat{\theta}$  has rank  $d$ . The first order condition  $\nabla \hat{L}(\theta; Z) = 0$  is the same as  $\nabla \hat{L}(g(\hat{\theta}); Z) = 0$ , which gives

$$\nabla \hat{L}(\hat{\phi}; Z) \hat{\theta} = 0. \quad (20)$$

As  $\hat{\theta}$  is of full row rank, there exists a matrix  $\hat{\theta}^\dagger \in \mathbb{R}^{k \times d}$  such that  $\hat{\theta} \hat{\theta}^\dagger = I$  (e.g.  $\hat{\theta}^\dagger = \hat{\theta}^\top (\hat{\theta} \hat{\theta}^\top)^{-1}$ ). We can right multiply the above equation by  $\hat{\theta}^\dagger$  and obtain that

$$\nabla \hat{L}(\hat{\phi}; Z) = 0.$$

As  $\hat{L}(\phi; Z)$  is convex in  $\phi$ , the above implies  $\hat{\phi} = g(\hat{\theta})$  is a global minimum of  $\hat{L}(\phi; Z)$ . Therefore,  $\hat{\theta}$  is a global minimum of  $\hat{L}(\theta; Z)$ . Next, consider the case where the rank of  $\hat{\theta}$  is smaller than  $d$ . In this case, we follow the proof strategy in Proposition 4 in [7]; we provide here for completeness. In this case, Equation (20) still holds, which implies

$$0 = \nabla \hat{L}(\hat{\phi}; Z) \hat{\theta} \hat{\theta}^\top = \nabla \hat{L}(\hat{\phi}; Z) \hat{\phi}. \quad (21)$$

The Hessian of  $\hat{L}(g(\hat{\theta}); Z)$  has

$$\nabla^2 \hat{L}(g(\hat{\theta}); Z)(\Delta, \Delta) = 2\langle \nabla \hat{L}(\hat{\phi}; Z), \Delta \Delta^\top \rangle + \nabla^2 \hat{L}(\hat{\phi}; Z)(\hat{\theta} \Delta^\top + \Delta \hat{\theta}^\top, \hat{\theta} \Delta^\top + \Delta \hat{\theta}^\top).$$

As  $\hat{\theta}R$  is also a local minimum for any orthogonal matrix  $R$  (i.e.,  $RR^\top = R^\top R = I$ ), we can find a  $\hat{\theta}$  with the last column being 0 by right multiplying certain  $R$ . Then, consider any  $\Delta$  with the first  $k-1$  columns being 0 and the last column being any  $u \in \mathbb{R}^d$ . With this choice of  $\Delta$  and  $\hat{\theta}$ ,  $\hat{\theta} \Delta^\top = 0$ . Therefore,

$$\nabla^2 \hat{L}(g(\hat{\theta}); Z)(\Delta, \Delta) = 2u^\top \nabla \hat{L}(\hat{\phi}; Z)u.$$

Since  $\hat{\theta}$  is a local minimum of  $\hat{L}(g(\hat{\theta}); Z)$ , it holds that  $\nabla^2 \hat{L}(g(\hat{\theta}); Z)(\Delta, \Delta) \geq 0$ , which implies

$$\nabla \hat{L}(\hat{\phi}; Z) \succeq 0. \quad (22)$$

Equation (21) and (22) together comprise the first order conditions of the convex minimization problem  $\min_{\phi \succeq 0} \hat{L}(\phi; Z)$ . Thus,  $\hat{\theta}$  is also a global minimum.  $\square$

## D Proofs for Section 4

First, we provide the full statement of Theorem 8 (including constants).

**Theorem 24.** *The expected regret of Algorithm 1 is*

$$R(T) \leq C_0 + C_1 \cdot T^{2/3} \left( \log \left( 3 + \frac{8\phi_{\max}KT}{\ell_{\max}^2} \right) \right)^{1/3},$$

where

$$C_0 = \frac{64(\phi_{\max})^{\frac{2d^2+2}{2d^2-1}}}{(135M(\ell_{\max} + \xi_{\max})^2d^3)^{\frac{2d^2+2}{2d^2-1}}(8\phi_{\max}K/\ell_{\max}^2)^{\frac{3d^2}{2d^2-1}}},$$

$$C_1 = 162d^2(M^2(\ell_{\max} + \xi_{\max})^4\phi_{\max})^{1/3}.$$

Before we prove Theorem 24, we first prove a preliminary result establishing an analog of the *smooth best arm response* property [39] to our setting. First, we have the following useful result:

**Lemma 25.** *Let  $\phi, \phi' \in \mathbb{R}^{d \times d}$  be symmetric matrices, let  $x, x' \in \mathbb{R}^d$  be eigenvectors of  $\phi, \phi'$  corresponding to their top eigenvalue, such that  $\|x\|_2 = \|x'\|_2 = 1$ , and let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$  be the eigenvalues of  $\phi'$ . Suppose that  $\langle x, x' \rangle \geq 0$ . Then, we have*

$$\|x - x'\|_2 \leq \frac{2^{3/2}\|\phi - \phi'\|_2}{\lambda_1 - \lambda_2}.$$

*Proof.* See Corollary 3 of [50]. □

Next, let  $\chi : \mathbb{R}^{d \times d} \rightarrow 2^{\mathcal{X}}$  denote the subset of reward-maximizing arms for  $g(\theta) = \theta\theta^\top$ —i.e.,

$$\chi(\phi) = \arg \max_{x \in \mathcal{X}} x^\top \phi x,$$

where the argmax returns the set of all optimal values. Then, we have the following analog of smooth best arm response:

**Lemma 26.** *For any  $\phi \in \mathbb{R}^{d \times d}$ , there exists  $x \in \chi(\phi)$  and  $x^* \in \chi(\phi^*)$  such that*

$$\|x - x^*\|_2 \leq M\|\phi - \phi^*\|_F.$$

*Proof.* First, note that  $x, x^*$  are eigenvectors of  $\phi, \phi^*$  corresponding to their top eigenvalues, respectively. Next, note that if  $x^* \in \chi(\phi^*)$ , then we also have  $-x^* \in \chi(\phi^*)$ ; thus, without loss of generality, we can assume that  $\langle x^*, x \rangle \geq 0$ . Also, note that  $\|x\|_2 = \|x^*\|_2 = 1$  since the optimizer maximizes the magnitude of  $x$ . Thus, we have

$$\|x - x^*\|_2 \leq \frac{2^{3/2}\|\phi - \phi^*\|_2}{\lambda_1 - \lambda_2} \leq M\|\phi - \phi^*\|_F,$$

where the second inequality follows by Lemma 25, and the third inequality follows by Assumption A3, as claimed. □



Now, we prove Theorem 24. The cumulative regret  $R(T)$  of a horizon of  $T$  has that

$$\begin{aligned}
R(T) &= \mathbb{E} \left[ \sum_{t=1}^T (f_{\theta^*}(x^*) - f_{\theta^*}(x_t)) \right] \\
&= \mathbb{E} \left[ \sum_{t=1}^m (f_{\theta^*}(x^*) - f_{\theta^*}(x_t)) + \sum_{t=m+1}^T (f_{\theta^*}(x^*) - f_{\theta^*}(x_t)) \right] \\
&\leq 2m\phi_{\max} + \mathbb{E} \left[ \sum_{t=m+1}^T \langle g(\hat{\theta}) - g(\theta^*), \hat{x}\hat{x}^\top - x^*x^{*\top} \rangle + \sum_{t=m+1}^T \langle g(\hat{\theta}), x^*x^{*\top} - \hat{x}\hat{x}^\top \rangle \right], \quad (23)
\end{aligned}$$

where  $\hat{\theta}$  is an estimator that minimizes the empirical loss of the first  $m$  samples,  $\hat{x} \in \chi(g(\hat{\theta}))$  maximizes the estimated expected reward  $f_{\hat{\theta}}(x)$ , and  $\langle \phi, \phi' \rangle = \sum_{i,j=1}^d \phi_{ij} \phi'_{ij}$  is the matrix inner product. Since  $\hat{x}$  is a maximizer of  $f_{\hat{\theta}}(x) = \langle g(\hat{\theta}), xx^\top \rangle$ , we have  $\langle g(\hat{\theta}), x^*x^{*\top} - \hat{x}\hat{x}^\top \rangle \leq 0$ . Thus, continuing from (23), we have

$$\begin{aligned}
R(T) &\leq 2m\phi_{\max} + \mathbb{E} \left[ \sum_{t=m+1}^T \langle g(\hat{\theta}) - g(\theta^*), \hat{x}\hat{x}^\top - x^*x^{*\top} \rangle \right] \\
&\leq 2m\phi_{\max} + (T - m)\mathbb{E} \left[ \|g(\hat{\theta}) - g(\theta^*)\|_F \|\hat{x}\hat{x}^\top - x^*x^{*\top}\|_F \right]. \quad (24)
\end{aligned}$$

To bound the second term in (24), note that

$$\|\hat{x}\hat{x}^\top - x^*x^{*\top}\|_F \leq \|\hat{x}\hat{x}^\top - \hat{x}x^{*\top}\|_F + \|\hat{x}x^{*\top} - x^*x^{*\top}\|_F \leq 2M\|g(\hat{\theta}) - g(\theta^*)\|_F,$$

where the last step follows by Lemma 26. Next, by Theorem 5, we have

$$\begin{aligned}
\|g(\hat{\theta}) - g(\theta^*)\|_F &\leq \sqrt{\frac{2\epsilon}{\alpha}} \\
&= d \left( \frac{45^3 \ell_{\max}^2 (\ell_{\max}^2 + \xi_{\max}^2)}{10m} \left( d^2 \max \left\{ 1, \log \left( 1 + \frac{8\phi_{\max} K m}{\ell_{\max}^2} \right) \right\} + \log \frac{2}{\delta} \right) \right)^{1/4}
\end{aligned}$$

with probability at least  $1 - \delta$ . Now, defining the event

$$\mathcal{G} = \left\{ \|g(\hat{\theta}) - g(\theta^*)\|_F \leq \sqrt{\frac{2\epsilon}{\alpha}} \right\},$$

letting

$$\delta = 2 \exp \left( -d^2 \max \left\{ 1, \log \left( 1 + \frac{8\phi_{\max} K m}{\ell_{\max}^2} \right) \right\} \right),$$

and continuing from (24), we have

$$\begin{aligned}
R(T) &\leq 2m\phi_{\max} + T \cdot \mathbb{E} \left[ \|g(\hat{\theta}) - g(\theta^*)\|_F \|\hat{x}\hat{x}^\top - x^*x^{*\top}\|_F \mathbb{1}(\mathcal{G}) \right] + 4T\phi_{\max} \cdot \mathbb{P}(\mathcal{G}^c) \\
&\leq 2m\phi_{\max} + 2MT \cdot \mathbb{E} \left[ \|g(\hat{\theta}) - g(\theta^*)\|_F^2 \mid \mathcal{G} \right] + 4T\phi_{\max} \cdot \mathbb{P}(\mathcal{G}^c) \\
&\leq 2m\phi_{\max} + 270M(\ell_{\max} + \xi_{\max})^2 d^3 T \sqrt{\frac{\log(3 + 8\phi_{\max} K T / \ell_{\max}^2)}{m}} + \frac{8T\phi_{\max}}{(8\phi_{\max} K m / \ell_{\max}^2)^{d^2}}. \quad (25)
\end{aligned}$$

The third term in inequality (25) is smaller than the second term when

$$m \geq \left( \frac{4\phi_{\max}}{135M(\ell_{\max} + \xi_{\max})^2 d^3 (8\phi_{\max}K/\ell_{\max}^2)^{d^2} \sqrt{\log(3 + 8\phi_{\max}KT/\ell_{\max}^2)}} \right)^{1/(d^2-1/2)}. \quad (26)$$

For a choice of  $m$  satisfying (26), continuing from (25), we have

$$R(T) \leq 2m\phi_{\max} + 540M(\ell_{\max} + \xi_{\max})^2 d^3 T \sqrt{\frac{\log(3 + 8\phi_{\max}KT/\ell_{\max}^2)}{m}}. \quad (27)$$

Next, we choose  $m$  to minimize the upper bound in (26) for sufficiently large  $T$ —in particular,

$$m = \left\lceil \left( \frac{135M(\ell_{\max} + \xi_{\max})^2 d^3 T \sqrt{\log(3 + 8\phi_{\max}KT/\ell_{\max}^2)}}{\phi_{\max}} \right)^{\frac{2}{3}} \right\rceil. \quad (28)$$

With this choice of  $m$ , we have

$$R(T) \leq 162(M^2(\ell_{\max} + \xi_{\max})^4 \phi_{\max})^{\frac{1}{3}} d^2 T^{2/3} (\log(3 + 8\phi_{\max}KT/\ell_{\max}^2))^{\frac{1}{3}}.$$

Finally, note that (26) holds under the choice of  $m$  in (28) for  $T$  satisfying

$$T \sqrt{\log(3 + 8\phi_{\max}KT/\ell_{\max}^2)} \geq \frac{64(\phi_{\max})^{\frac{2d^2+2}{2d^2-1}}}{(135M(\ell_{\max} + \xi_{\max})^2 d^3)^{\frac{2d^2+2}{2d^2-1}} (8\phi_{\max}K/\ell_{\max}^2)^{\frac{3d^2}{2d^2-1}}}.$$

The claim follows.  $\square$

## E Proof of Theorem 9

First, we have the following key result:

**Lemma 27.** *Let  $\theta, \theta' \in \mathbb{R}^{d \times k}$ , and let  $\phi = \theta\theta^\top$  and  $\phi' = \theta'\theta'^\top$ . Assume that  $\|\phi - \phi'\|_F \leq \eta$ , and that  $\sigma_{\min}(\theta) \geq \sigma_0 > 0$ , where  $\sigma_{\min}(\theta)$  is the minimum singular value of  $\theta$  (more precisely, the  $d$ th largest singular value). Then, there exist orthogonal matrices  $R, R' \in \mathbb{R}^{k \times k}$  such that*

$$\|\theta R - \theta' R'\|_F \leq \frac{\eta}{\sigma_0}. \quad (29)$$

*Proof.* Consider the SVDs  $\theta = U\Sigma V^\top$  and  $\theta' = U'\Sigma'V'^\top$ , where  $U, U' \in \mathbb{R}^{d \times d}$ ,  $\Sigma, \Sigma' \in \mathbb{R}^{d \times d}$ , and  $V, V' \in \mathbb{R}^{k \times d}$ ; then, we have  $\phi = U\Sigma^2 U^\top$  and  $\phi' = U'\Sigma'^2 U'^\top$ . Then, we claim that the choices  $R = VU^\top$  and  $R' = V'U'^\top$  satisfy (29). In particular, note that  $\theta R = U\Sigma U^\top$  and  $\theta' R' = U'\Sigma'U'^\top$ , since  $V^\top V = V'^\top V' = I_d$  since  $k \geq d$ , where  $I_d \in \mathbb{R}^{d \times d}$  is the  $d$ -dimensional identity matrix. Thus, it suffices to show that

$$\sigma_0 \|U\Sigma U^\top - U'\Sigma'U'^\top\|_F \leq \eta. \quad (30)$$

To this end, note that

$$\eta \geq \|\phi - \phi'\|_F = \|U\Sigma^2 U^\top - U'\Sigma'^2 U'^\top\|_F = \|U'^\top U\Sigma^2 - \Sigma'^2 U'^\top U\|_F, \quad (31)$$

where in the last step, we have multiplied the expression inside the Frobenius norm by  $U'^\top$  on the left and by  $U$  on the right, using the fact that the Frobenius norm is invariant under multiplication by orthogonal matrices. Defining  $W = U'^\top U$ , note that

$$(W\Sigma)_{ij} = \sum_{k=1}^d W_{ik}\Sigma_{kj} = W_{ij}\Sigma_{jj} \quad (32)$$

$$(\Sigma W)_{ij} = \sum_{k=1}^d \Sigma'_{ik}W_{kj} = W_{ij}\Sigma'_{ii} \quad (33)$$

$$(W\Sigma^2)_{ij} = \sum_{k=1}^d W_{ik}(\Sigma^2)_{kj} = W_{ij}\Sigma_{jj}^2 \quad (34)$$

$$(\Sigma^2 W)_{ij} = \sum_{k=1}^d (\Sigma'^2)_{ik}W_{kj} = W_{ij}\Sigma'_{ii}{}^2. \quad (35)$$

Then, continuing from (31), we have

$$\begin{aligned} \eta^2 &\geq \|W\Sigma^2 - \Sigma'^2 W\|_F^2 = \sum_{i,j=1}^d W_{ij}^2 (\Sigma_{jj}^2 - \Sigma'_{ii}{}^2)^2 \\ &= \sum_{i,j=1}^d W_{ij}^2 (\Sigma_{jj} - \Sigma'_{ii})^2 (\Sigma_{jj} + \Sigma'_{ii})^2 \\ &\geq \sum_{i,j=1}^d W_{ij}^2 (\Sigma_{jj} - \Sigma'_{ii})^2 \sigma_0^2 \\ &= \sigma_0^2 \|W\Sigma - \Sigma' W\|_F^2 \\ &= \sigma_0^2 \|U'^\top U\Sigma - \Sigma' U'^\top U\|_F^2 \\ &= \sigma_0^2 \|U\Sigma U^\top - U'\Sigma' U'^\top\|_F^2, \end{aligned}$$

where on the first line, we have used (34) & (35), on the third line we have used  $\Sigma_{jj} \geq \sigma_0$ , on the fourth line we have used (32) & (33), and on the last line we have multiplied on by  $U'$  on the left  $U'^\top$  on the right, again using the fact that the Frobenius norm is invariant under multiplication by orthogonal matrices. Thus, we have shown (30), so the claim follows.  $\square$

We note here that our result provides an analog of Lemma 6 in [20] for quadratic neural networks.

Now, we prove Theorem 9. First, by directly applying the arguments in the proof of Theorem 5, we have

$$\|g(\hat{\theta}_p) - g(\theta_p^*)\|_F \leq \sqrt{\frac{2\epsilon_p}{\alpha}}$$

with probability at least  $1 - \delta/2$ . However,  $\hat{\theta}_p$  itself may not be close to  $\theta_p^*$ . Instead, applying Lemma 27 with  $\theta = \hat{\theta}_p$  and  $\theta' = \theta_p^*$ , and with  $\eta = \sqrt{2\epsilon_p/\alpha}$ , there exists an orthogonal matrix  $R_p = R'R^\top$  that ‘‘aligns’’  $\hat{\theta}_p$  with  $\theta_p^*$ , yielding

$$\|\hat{\theta}_p - \theta_p^* R_p\|_F \leq \frac{1}{\sigma_0} \sqrt{\frac{2\epsilon_p}{\alpha}},$$

where  $\sigma_0$  is the minimum singular value of  $\theta_p^*$ . Now, let  $\tilde{\theta}_g = \theta_g^* R_p$ , and note that this is a global minimizer (i.e.,  $g(\tilde{\theta}_g) = g(\theta_g^*)$ ), since  $R_p$  is orthogonal. Then, we have

$$\begin{aligned} \|\tilde{\theta}_g - \hat{\theta}_p\|_F &\leq \|\theta_g^* R_p - \theta_p^* R_p\|_F + \|\theta_p^* R_p - \hat{\theta}_p\|_F \\ &\leq \|\theta_g^* - \theta_p^*\|_F + \frac{1}{\sigma_0} \sqrt{\frac{2\epsilon_p}{\alpha}} \\ &\leq B + \frac{1}{\sigma_0} \sqrt{\frac{2\epsilon_p}{\alpha}} \end{aligned} \quad (36)$$

with probability at least  $1 - \delta/2$ . In other words, an alternative global minimizer  $\tilde{\theta}_g$  exists within a small Frobenius norm of our proxy estimator  $\hat{\theta}_p$ , even if  $\hat{\theta}_p$  is not close to  $\theta_p^*$ .

Finally, on the event that (36) holds, note that for  $\theta \in B_2(\hat{\theta}_p, \hat{B})$ , we have the alternative upper bound

$$|f_\theta(x) - f_{\theta_g^*}(x)| \leq K \|g(\theta) - g(\theta_g^*)\|_F \leq K \hat{B},$$

where the first inequality holds by Lemma 3; thus, we can take  $\ell_{\max} = K \hat{B}$ . Thus, on the event that (36) holds, by Theorem 5, we have

$$\mathbb{P}_{p(Z)} \left[ L_q(\hat{\theta}_g) \leq \frac{2K^2 \epsilon_g}{\alpha} \right] \geq 1 - \frac{\delta}{2},$$

so the claim follows by a union bound.  $\square$

## F Proofs for Section 6

### F.1 Proof of Lemma 13

We prove the case

$$\int_{X_i} |f_\theta(x) - f_{\theta^*}(x)| dx - \int_{X'_i} |f_\theta(x) - f_{\theta^*}(x)| dx \leq 2k\epsilon^3 \sqrt{d} \cdot |S^{d-2}|.$$

the proof of the negation is identical. First, note that

$$\begin{aligned} \int_{X'_i} |f_\theta(x) - f_{\theta^*}(x)| dx &= \int_{X_i} |f_\theta(\phi(x)) - f_{\theta^*}(\phi(x))| \cdot |\det \nabla_x \phi(x)| dx \\ &= \int_{X_i} |f_\theta(\phi(x)) - f_{\theta^*}(\phi(x))| dx, \end{aligned}$$

since  $\nabla_x \phi(x)$  is a lower triangular matrix with all ones along its diagonal. Now, note that

$$\begin{aligned} &\int_{X_i} |f_\theta(x) - f_{\theta^*}(x)| - |f_\theta(\phi(x)) - f_{\theta^*}(\phi(x))| dx \\ &\leq \int_{X_i} |(f_\theta(x) - f_\theta(\phi(x))) - (f_{\theta^*}(x) - f_{\theta^*}(\phi(x)))| dx \\ &\leq \int_{X_i} |f_\theta(x) - f_\theta(\phi(x))| + |f_{\theta^*}(x) - f_{\theta^*}(\phi(x))| dx \\ &\leq 2L |X_i| \cdot \max_{x \in X_i} \|x - \phi(x)\|_2, \end{aligned}$$

where  $L$  is a Lipschitz constant for  $f_\theta$  as a function of  $x$ . Next, note that

$$\|x - \phi(x)\|_2 \leq \left( \frac{1}{\sqrt{1-x_1^2}} - 1 \right) \sqrt{d} \leq \frac{\epsilon^2 \sqrt{d}}{2}$$

for all  $x \in X_i$ . Finally, note that for any  $x, x' \in \mathbb{R}^d$ , we have

$$|f_\theta(x) - f_\theta(x')| = \left| \sum_{i=1}^k \sigma(\theta_i^\top x) - \sigma(\theta_i^\top x') \right| \leq \sum_{i=1}^k |\theta_i^\top (x - x')| \leq k \|x - x'\|_2,$$

which implies that  $L \leq k$ . Finally, note that

$$|X_i| = \int_{-\epsilon}^{\epsilon} V^{d-2} \left( \sqrt{1-x_1^2} \right) dx_1 \leq 2\epsilon \cdot |S^{d-2}|,$$

where  $V^n(r)$  is the volume of the  $n$ -sphere with radius  $r$ , so the claim follows.  $\square$

## F.2 Proof of Lemma 14

Note that

$$\begin{aligned} \mathcal{F}(a, b) &= \int_0^\epsilon |a - bw| dw \\ &= \int_0^{a/b} (a - bw) dw + \int_{a/b}^\epsilon (bw - a) dw \\ &= \left[ aw - \frac{bw^2}{2} \right]_0^{a/b} + \left[ \frac{bw^2}{2} - aw \right]_{a/b}^\epsilon \\ &= \left( \frac{a^2}{b} - \frac{a^2}{2b} \right) + \left( \frac{b\epsilon^2}{2} - a\epsilon \right) - \left( \frac{a^2}{2b} - \frac{a^2}{b} \right) \\ &= \frac{a^2}{b} + \frac{b\epsilon^2}{2} - a\epsilon. \end{aligned}$$

As a function of  $a$ , this expression is minimized when  $a = b\epsilon/2$ , in which case

$$\mathcal{F}\left(\frac{b\epsilon}{2}, b\right) = \frac{b\epsilon^2}{2}.$$

Now, note that

$$\begin{aligned} \int_{-\epsilon}^\epsilon |(\beta_0 + \beta_1 w) - g(w)| dw &= \int_{-\epsilon}^0 |\beta_0 + \beta_1 w| dw + \int_0^\epsilon |\beta_0 + (\beta_1 - 1)w| dw \\ &= \int_0^\epsilon |\beta_0 - \beta_1 w| dw + \int_0^\epsilon |\beta_0 + (\beta_1 - 1)w| dw \\ &= \mathcal{F}(\beta_0, \beta_1) + \mathcal{F}(\beta_0, 1 - \beta_1) \end{aligned}$$

Now, we must have either  $\beta_1 \geq 1/2$  or  $1 - \beta_1 \geq 1/2$ ; without loss of generality, assume the former holds. Then, we have

$$\int_{-\epsilon}^\epsilon |(\beta_0 + \beta_1 w) - g(w)| dw \geq \mathcal{F}(\beta_0, \beta_1) \geq \frac{\beta_1 \epsilon^2}{2} \geq \frac{\epsilon^2}{4},$$

as claimed.  $\square$

### F.3 Proof of Lemma 15

Without loss of generality, we can modify the coordinate system so that

$$\beta = [t_1 \ t_2 \ 0 \ \dots \ 0]^\top$$

without affecting  $\theta_i^*$ . By assumption, we have  $\|\beta\|_2 = \sqrt{t_1^2 + t_2^2} = 1$ . In the following argument, we consider the case  $t_1 \geq 0$ ; the case  $t_1 < 0$  follows from the same argument with  $\alpha^2 < \|\beta + \theta_i^*\|_2^2$ . Now, we have

$$\begin{aligned} \alpha^2 < \|\beta - \theta_i^*\|_2^2 &= (1 - t_1)^2 + t_2^2 \\ &= 2(1 - t_1) \\ &= 2 \left( 1 - \frac{t_1}{\sqrt{t_1^2 + t_2^2}} \right) \\ &= 2 \left( 1 - \frac{1}{\sqrt{1 + t_2^2/t_1^2}} \right), \end{aligned}$$

so

$$\frac{|t_2|}{|t_1|} > \sqrt{\left(\frac{1}{1 - \alpha^2/2}\right)^2 - 1} \geq \sqrt{\left(1 + \frac{\alpha^2}{2}\right)^2 - 1} \geq \alpha$$

Next, the condition  $\beta^\top([x_1] \circ z) = 0$  is equivalent to

$$t_1 x_1 + t_2 z_1 = 0,$$

or

$$z_1 = -\frac{t_1 x_1}{t_2}.$$

As a consequence, we have

$$|z_1| \leq \frac{|t_1| \cdot |x_1|}{|t_2|} \leq \frac{\epsilon}{\alpha}.$$

Thus, letting

$$V^{d-3}(r) = \text{Vol}(\{w \in \mathbb{R}^{d-2} \mid \|w\|_2 = r\})$$

be the volume of the  $d - 3$  sphere of radius  $r$ , then we have

$$|Z_i^\beta| \leq \int_{-\epsilon/\alpha}^{\epsilon/\alpha} V^{d-3} \left( \sqrt{1 - z_1^2} \right) dz_1 \leq \int_{-\epsilon/\alpha}^{\epsilon/\alpha} V^{d-3}(1) dz_1 = \frac{2\epsilon \cdot |S^{d-3}|}{\alpha},$$

as claimed.  $\square$

#### F.4 Proof of Theorem 11

**Lemma 28.** *Under Assumptions A5,  $f_\theta$  and  $L$  are  $K$ -Lipschitz in  $\theta$  with respect to  $\ell_{2,1}$  norm, where  $K = 4k$  and the  $\ell_{2,1}$  norm for any parameter  $\theta \in \mathbb{R}^{d \times k}$  is defined as  $\sum_{i=1}^k \|\theta_i\|$ .*

*Proof.* By our definition, for any  $\theta, \theta' \in \Theta$ ,

$$|f_\theta(x) - f_{\theta'}(x)| = \left| \sum_{i=1}^k \sigma(\theta_i^\top x) - \sum_{i=1}^k \sigma(\theta'_i{}^\top x) \right| \leq \sum_{i=1}^k \|\theta_i - \theta'_i\|_2.$$

Given our quadratic loss function, we have

$$\begin{aligned} & |(f_\theta(x) - f_{\theta^*}(x))^2 - (f_{\theta'}(x) - f_{\theta^*}(x))^2| \\ & \leq |f_\theta(x) - f_{\theta^*}(x) + f_{\theta'}(x) - f_{\theta^*}(x)| |f_\theta(x) - f_{\theta'}(x)| \\ & \leq 4k \sum_{i=1}^k \|\theta_i - \theta'_i\|_2. \end{aligned}$$

Next, the true loss satisfies

$$|L_p(\theta) - L_p(\theta')| \leq \mathbb{E}_{p(x)}[|(f_\theta(x) - f_{\theta^*}(x))^2 - (f_{\theta'}(x) - f_{\theta^*}(x))^2|] \leq 4k \sum_{i=1}^k \|\theta_i - \theta'_i\|_2.$$

Finally, the empirical loss satisfies

$$\begin{aligned} |\hat{L}(\theta; Z) - \hat{L}(\theta'; Z)| &= \left| \frac{1}{n} \sum_{i=1}^n [(f_\theta(x_i) - y_i)^2 - (f_{\theta'}(x_i) - y_i)^2] \right| \\ &\leq \frac{1}{n} \sum_{i=1}^n |(f_\theta(x_i) - f_{\theta^*}(x_i))^2 - (f_{\theta'}(x_i) - f_{\theta^*}(x_i))^2| + \frac{2}{n} \sum_{i=1}^n |\xi_i| \cdot |f_\theta(x_i) - f_{\theta'}(x_i)| \\ &\leq (4k + 2\xi_{\max}) \sum_{i=1}^k \|\theta_i - \theta'_i\|_2, \end{aligned}$$

as claimed. □

**Lemma 29.** *Under Assumptions A5, for any  $\delta \in \mathbb{R}_{>0}$ , we have*

$$\mathbb{P}_{p(Z)} \left[ \sup_{\theta \in \Theta} |L_p(\theta) - \hat{L}(\theta; Z) - \sigma(Z)| \leq \epsilon \right] \geq 1 - \delta,$$

where  $\sigma(Z) = n^{-1} \sum_{i=1}^n \xi_i^2$ , and letting  $\ell_{\max} = 2k$  be an upper bound on  $|f_\theta(x) - f_{\theta^*}(x)|$ , where  $\epsilon$  is as in Theorem 11.

*Proof.* Consider an  $\epsilon/(4K)$ -net  $\mathcal{E}$  with respect to  $\ell_{2,1}$  norm. Then, for any  $\theta \in \Theta$ , there exists  $\theta' \in \mathcal{E}$  such that

$$|(\hat{L}(\theta; Z) - L_p(\theta)) - (\hat{L}(\theta'; Z) - L_p(\theta'))| \leq 2K \sum_{i=1}^k \|\theta_i - \theta'_i\|_2 \leq \frac{\epsilon}{2}.$$

Therefore, we have

$$\begin{aligned}
& \mathbb{P}_{p(Z)} \left[ \sup_{\theta} |\hat{L}(\theta; Z) - L_p(\theta) - \sigma(Z)| \geq \epsilon \right] \\
& \leq \mathbb{P}_{p(Z)} \left[ \max_{\theta \in \mathcal{E}} |\hat{L}(\theta; Z) - L_p(\theta) - \sigma(Z)| \geq \frac{\epsilon}{2} \right] \\
& \leq \sum_{\theta \in \mathcal{E}} \mathbb{P}_{p(Z)} \left[ |\hat{L}(\theta; Z) - L_p(\theta) - \sigma(Z)| \geq \frac{\epsilon}{2} \right].
\end{aligned} \tag{37}$$

Following a similar argument as in Lemma 4, we obtain from (37) that

$$\begin{aligned}
& \sum_{\theta \in \mathcal{E}} \mathbb{P}_{p(Z)} \left[ |\hat{L}(\theta; Z) - L_p(\theta) - \sigma(Z)| \geq \frac{\epsilon}{2} \right] \\
& \leq 4|\mathcal{E}| \cdot \exp \left( -\frac{n\epsilon^2}{18\ell_{\max}^2(\ell_{\max}^2 + \xi_{\max}^2)} \right) \\
& \leq 2 \left( 1 + \frac{4kK}{\epsilon} \right)^{dk} \cdot \exp \left( -\frac{n\epsilon^2}{18\ell_{\max}^2(\ell_{\max}^2 + \xi_{\max}^2)} \right) \\
& = 2 \exp \left( -\frac{n\epsilon^2}{18\ell_{\max}^2(\ell_{\max}^2 + \xi_{\max}^2)} + dk \log \left( 1 + \frac{4kK}{\epsilon} \right) \right),
\end{aligned} \tag{38}$$

where the second inequality follows since by Lemma 22, the covering number of the  $\epsilon$ -net  $\mathcal{E}$  of  $\Theta$  satisfies

$$|\mathcal{E}| \leq \left( 1 + \frac{k}{\epsilon} \right)^{dk}.$$

Finally, we choose  $\epsilon$  so that (38) is smaller than  $\delta$ —in particular, letting

$$\epsilon = \sqrt{\frac{18\ell_{\max}^2(\ell_{\max}^2 + \xi_{\max}^2)}{n} \left( dk \max \left\{ 1, \log \left( 1 + \frac{4kKn}{\ell_{\max}^2} \right) \right\} + \log \frac{2}{\delta} \right)}.$$

then continuing (38), we have

$$2 \exp \left( -\frac{n\epsilon^2}{18\ell_{\max}^2(\ell_{\max}^2 + \xi_{\max}^2)} + dk \log \left( 1 + \frac{4kK}{\epsilon} \right) \right) \leq \delta,$$

as claimed.  $\square$

$\square$

Finally, to prove Theorem 11, note that

$$L_p(\hat{\theta}(Z)) \leq \hat{L}(\theta^*; Z) + \epsilon - \sigma(Z) \leq L_p(\theta^*) + 2\epsilon = 2\epsilon$$

with probability at least  $1 - \delta$ . Thus,

$$\mathbb{E}_{p(x)} [|f_{\hat{\theta}}(x) - f_{\theta^*}(x)|] \leq (\mathbb{E}_{p(x)} [(f_{\hat{\theta}}(x) - f_{\theta^*}(x))^2])^{\frac{1}{2}} = L_p(\hat{\theta}(Z)) \leq \sqrt{2\epsilon}.$$

Then by Lemma 10, we have for any  $x \in \mathcal{X}$

$$|f_{\hat{\theta}}(x) - f_{\theta^*}(x)| \leq 20k^2 \sqrt{d^3(2\epsilon)^{1/2}},$$



that is,

$$\mathbb{P}_{p(Z)} \left[ L_q(\hat{\theta}(Z)) \leq 400k^4 d^3 (2\epsilon)^{1/2} \right] \geq 1 - \delta.$$

Finally, we note that to satisfy the condition  $\eta \leq (6126d^2k^2)^{-1}$ , it suffices to have

$$n \geq 72 \cdot 6126^4 \ell_{\max}^2 (\ell_{\max}^2 + \xi_{\max}^2) d^8 k^8 \left( dk \max \left\{ 1, \log \left( 1 + \frac{4kKn}{\ell_{\max}^2} \right) \right\} + \log \frac{2}{\delta} \right).$$

The claim follows.  $\square$

## G Proofs for Appendix B

### G.1 Proof of Proposition 18

Suppose that  $z_t \in \{0, 1\}$  is binary,  $z_0 = 0$ , and

$$\tilde{p}(z | z') = \begin{cases} 1 & \text{if } z = z' \\ 0 & \text{otherwise.} \end{cases}$$

In particular, since  $z_0 = 0$ ,  $p(w) = \mathbb{1}(w = w_0)$  places all weight on the zero sequence  $w_0 = 0\dots 0$ . Next, consider the shifted distribution

$$\tilde{q}(z_t | z_{t-1}) = \begin{cases} 1 & \text{if } z = z' = 1 \\ 1 - \alpha/2 & \text{if } z = z' = 0 \\ \alpha/2 & \text{otherwise.} \end{cases}$$

Note that  $\|\tilde{p}(\cdot | z') - \tilde{q}(\cdot | z')\|_{\text{TV}} \leq \alpha$ , so Assumption A7 is satisfied. Note that

$$q(w_0) = \prod_{t=1}^T \tilde{q}(0 | 0) = (1 - \alpha/2)^T.$$

As a consequence, we have

$$\begin{aligned} \|p - q\|_{\text{TV}} &= \sum_{w \in \mathcal{W}} |p(w) - q(w)| \\ &= |p(w_0) - q(w_0)| + \sum_{w \in \mathcal{W} \setminus \{w_0\}} q(w) \\ &= (1 - (1 - \alpha/2)^T) + (1 - (1 - \alpha/2)^T) \\ &= 2(1 - (1 - \alpha/2)^T), \end{aligned}$$

as claimed.  $\square$

## G.2 Proof of Lemma 19

Note that

$$\begin{aligned}
& \|\tilde{q}_t - \tilde{p}_t\|_{\text{TV}} \\
&= \sum_{j=1}^k \int |\tilde{q}_t(z, j) - \tilde{p}_t(z, j)| dz \\
&= \sum_{j=1}^k \sum_{j'=1}^k \int \mathbb{1}(j = \tilde{g}^*(z', j')) \cdot |\tilde{q}(z | z') \tilde{q}_{t-1}(z', j') - \tilde{p}(z | z') \tilde{p}_{t-1}(z', j')| dz' dz \\
&= \sum_{j'=1}^k \int |\tilde{q}(z | z') \tilde{q}_{t-1}(z', j') - \tilde{p}(z | z') \tilde{p}_{t-1}(z', j')| dz' dz \\
&\leq \sum_{j'=1}^k \int |\tilde{q}(z | z') - \tilde{p}(z | z')| \cdot \tilde{q}_{t-1}(z', j') + \tilde{p}(z | z') \cdot |\tilde{q}_{t-1}(z', j') - \tilde{p}_{t-1}(z', j')| dz' dz \\
&\leq \sum_{j'=1}^k \int \alpha \cdot \tilde{q}_{t-1}(z', j') + |\tilde{q}_{t-1}(z', j') - \tilde{p}_{t-1}(z', j')| dz' \\
&\leq \alpha + \|\tilde{q}_{t-1} - \tilde{p}_{t-1}\|_{\text{TV}}.
\end{aligned}$$

Since  $q_0(z, j) = p_0(z, j)$  for all  $z \in \mathcal{Z}$  and  $j \in [k]$ , by induction,  $\|\tilde{q}_t - \tilde{p}_t\|_{\text{TV}} \leq t\alpha$ . Thus, we have

$$\|\tilde{q} - \tilde{p}\|_{\text{TV}} \leq \frac{1}{T} \sum_{t=1}^T \|\tilde{q}_t - \tilde{p}_t\| \leq T\alpha,$$

as claimed.  $\square$

## G.3 Proof of Lemma 20

First, we prove the following lemma.

**Lemma 30.** *We have*

$$\tilde{p}_t(z_t, j_{t-1}) = \sum_{j_1, \dots, j_{t-2}}^k \int \left( \prod_{\tau=1}^{t-1} \mathbb{1}(j_\tau = \tilde{g}^*(z_\tau, j_{\tau-1})) \right) \cdot p(z_1, \dots, z_t) dz_1 \dots dz_{t-1}.$$

*Proof.* For the base case, we have

$$\begin{aligned}
\tilde{p}_2(z_2, j_1) &= \sum_{j_0=1}^k \int \mathbb{1}(j_1 = \tilde{g}^*(z_1, j_0)) \cdot \tilde{p}(z_2 | z_1) \cdot \tilde{p}_1(z_1, j_0) dz_1 \\
&= \sum_{j_0=1}^k \int \mathbb{1}(j_1 = \tilde{g}^*(z_1, j_0)) \cdot \tilde{p}(z_2 | z_1) \cdot \mathbb{1}(j = 0) \cdot \tilde{p}(z_1) dz_1 \\
&= \int \mathbb{1}(j_1 = \tilde{g}^*(z_1, j_0)) \cdot p(z_1, z_2) dz_1,
\end{aligned}$$

as claimed. For the inductive case, we have

$$\begin{aligned}\tilde{p}_t(z_t, j_{t-1}) &= \sum_{j_{t-2}=1}^k \int \mathbb{1}(j_{t-1} = \tilde{g}^*(z_{t-1}, j_{t-2})) \cdot \tilde{p}(z_t | z_{t-1}) \cdot \tilde{p}_{t-1}(z_{t-1}, j_{t-2}) dz_{t-1} \\ &= \sum_{j_1, \dots, j_{t-2}=1}^k \int \left( \prod_{\tau=1}^{t-1} \mathbb{1}(j_{\tau-1} = \tilde{g}^*(z_{\tau-1}, j_{\tau-2})) \right) \cdot p(z_1, \dots, z_t) dz_1 \dots dz_{t-1}.\end{aligned}$$

as claimed.  $\square$

Now, we prove Lemma 20. First, note that for each  $t \in [T]$ , we have

$$\begin{aligned}\mathbb{P}_{p(w)} \left[ (\hat{g}(w)_t \neq g^*(w)_t) \wedge \left( \bigwedge_{\tau=1}^{t-1} \hat{g}(w)_\tau = g^*(w)_\tau \right) \right] \\ &= \int \mathbb{1}(\hat{g}(w)_t \neq g^*(w)_t) \cdot \left( \prod_{\tau=1}^{t-1} \mathbb{1}(\hat{g}(w)_\tau = g^*(w)_\tau) \right) \cdot p(w) dw \\ &= \sum_{j_1, \dots, j_{t-1}=1}^k \int \mathbb{1}(\hat{g}(w)_t \neq g^*(w)_t) \cdot \left( \prod_{\tau=1}^{t-1} \mathbb{1}(\hat{g}(w)_\tau = g^*(w)_\tau) \right) \cdot p(j_1 \dots j_{t-1} | w) \cdot p(w) dw \\ &= \sum_{j_1, \dots, j_{t-1}=1}^k \int \mathbb{1}(\hat{g}(w)_t \neq g^*(w)_t) \cdot \left( \prod_{\tau=1}^{t-1} \mathbb{1}(\hat{g}(w)_\tau = g^*(w)_\tau) \right) \cdot \left( \prod_{\tau=1}^{t-1} \mathbb{1}(j_\tau = \tilde{g}^*(z_\tau, j_{\tau-1})) \right) \\ &\quad \cdot p(w) dw \\ &= \sum_{j_1, \dots, j_{t-1}=1}^k \int \mathbb{1}(\hat{g}(z_t, j_{t-1}) \neq \tilde{g}^*(z_t, j_{t-1})) \cdot \left( \prod_{\tau=1}^{t-1} \mathbb{1}(\hat{g}(z_\tau, j_{\tau-1}) = \tilde{g}^*(z_\tau, j_{\tau-1})) \right) \\ &\quad \cdot \left( \prod_{\tau=1}^{t-1} \mathbb{1}(j_\tau = \tilde{g}^*(z_\tau, j_{\tau-1})) \right) \cdot p(w) dw \\ &\leq \sum_{j_1, \dots, j_{t-1}=1}^k \int \mathbb{1}(\hat{g}(z_t, j_{t-1}) \neq \tilde{g}^*(z_t, j_{t-1})) \cdot \left( \prod_{\tau=1}^{t-1} \mathbb{1}(j_\tau = \tilde{g}^*(z_\tau, j_{\tau-1})) \right) \cdot p(w) dw \\ &= \sum_{j_1, \dots, j_{t-1}=1}^k \int \mathbb{1}(\hat{g}(z_t, j_{t-1}) \neq \tilde{g}^*(z_t, j_{t-1})) \cdot \left( \prod_{\tau=1}^{t-1} \mathbb{1}(j_\tau = \tilde{g}^*(z_\tau, j_{\tau-1})) \right) \cdot p(z_1, \dots, z_t) dz_1 \dots dz_t \\ &= \mathbb{P}_{p_t(z, j)} \left[ \hat{g}(z_t, j_{t-1}) \neq \tilde{g}^*(z_t, j_{t-1}) \right],\end{aligned}$$

where the last step follows from Lemma 30. Now, note that

$$\begin{aligned}\mathbb{P}_{p(w)}[\hat{g}(w) \neq g^*(w)] &= \sum_{t=1}^T \mathbb{P}_{p(w)} \left[ (\hat{g}(w)_t \neq g^*(w)_t) \wedge \left( \bigwedge_{\tau=1}^{t-1} \hat{g}(w)_\tau = g^*(w)_\tau \right) \right] \\ &\leq \sum_{t=1}^T \mathbb{P}_{\tilde{p}_t(z, j)} \left[ \hat{g}(z, j) \neq \tilde{g}^*(z, j) \right] \\ &\leq T \epsilon_g,\end{aligned}$$

as claimed.  $\square$

## G.4 Proof of Theorem 21

First, we show that  $\mathbb{P}_{q(z,j)}[\hat{g}(z,j) \neq \tilde{g}^*(z,j)] \leq \epsilon_g + T\alpha$ . To this end, note that

$$\begin{aligned}
& \mathbb{P}_{q(z,j)}[\hat{g}(z,j) \neq \tilde{g}^*(z,j)] \\
&= \mathbb{P}_{p(z,j)}[\hat{g}(z,j) \neq \tilde{g}^*(z,j)] + \mathbb{P}_{q(z,j)}[\hat{g}(z,j) \neq \tilde{g}^*(z,j)] - \mathbb{P}_{p(z,j)}[\hat{g}(z,j) \neq \tilde{g}^*(z,j)] \\
&\leq \epsilon_g + \sum_{j=1}^k \int \mathbb{1}(\hat{g}(z,j) \neq \tilde{g}^*(z,j)) \cdot |\tilde{q}(z,j) - \tilde{p}(z,j)| dz \\
&\leq \epsilon_g + \|\tilde{q} - \tilde{p}\|_{\text{TV}} \\
&\leq \epsilon_g + T\alpha.
\end{aligned}$$

Next, by Lemma 20 with  $q$  in place of  $p$  and  $\epsilon_g + T\alpha$  in place of  $\epsilon_g$ , we have  $\mathbb{P}_{q(w)}[\hat{g}(w) \neq g^*(w)] \leq T\epsilon_g + T^2\alpha$ . Then, assuming that  $\hat{g}(w) = g^*(w)$ , we have

$$\begin{aligned}
& \|f^*(x,w) - \hat{f}(x,w)\|_2 \\
&= \|(f_{j_T}^* \circ \dots \circ f_{j_1}^*)(x) - (\hat{f}_{j_T} \circ \dots \circ \hat{f}_{j_1})(x)\|_2 \\
&\leq \sum_{t=1}^T \|(f_{j_T}^* \circ \dots \circ f_{j_{t+1}}^* \circ f_{j_t}^* \circ \hat{f}_{j_{t-1}} \circ \dots \circ \hat{f}_{j_1})(x) - (f_{j_T}^* \circ \dots \circ f_{j_{t+1}}^* \circ \hat{f}_{j_t} \circ \hat{f}_{j_{t-1}} \circ \dots \circ \hat{f}_{j_1})(x)\|_2 \\
&\leq \sum_{t=1}^T K^{T-t} \cdot \|(f_{j_t}^* \circ \hat{f}_{j_{t-1}} \circ \dots \circ \hat{f}_{j_1})(x) - (\hat{f}_{j_t} \circ \hat{f}_{j_{t-1}} \circ \dots \circ \hat{f}_{j_1})(x)\|_2 \\
&\leq \sum_{t=1}^T K^{T-t} \epsilon_f \\
&\leq T\epsilon_f \cdot \max\{K^{T-1}, 1\}.
\end{aligned}$$

The claim follows by a union bound.  $\square$