# Gradient Descent with Early Stopping is Provably Robust to Label Noise for Overparameterized Neural Networks

Mingchen Li\* Mahdi Soltanolkotabi<sup>†</sup> Samet Oymak<sup>‡</sup>
March 31, 2019

#### Abstract

Modern neural networks are typically trained in an over-parameterized regime where the parameters of the model far exceed the size of the training data. Due to over-parameterization these neural networks in principle have the capacity to (over)fit any set of labels including pure noise. Despite this high fitting capacity, somewhat paradoxically, neural network models trained via first-order methods continue to predict well on yet unseen test data. In this paper we take a step towards demystifying this phenomena. In particular we show that first order methods such as gradient descent are provably robust to noise/corruption on a constant fraction of the labels despite over-parametrization under a rich dataset model. In particular: i) First, we show that in the first few iterations where the updates are still in the vicinity of the initialization these algorithms only fit to the correct labels essentially ignoring the noisy labels. ii) Secondly, we prove that to start to overfit to the noisy labels these algorithms must stray rather far from from the initial model which can only occur after many more iterations. Together, these show that gradient descent with early stopping is provably robust to label noise and shed light on empirical robustness of deep networks as well as commonly adopted heuristics to prevent overfitting.

## 1 Introduction

#### 1.1 Motivation

Deep neural networks (DNN) are ubiquitous in a growing number of domains ranging from computer vision to healthcare. State-of-the-art DNN models are typically overparameterized and contain more parameters than the size of the training dataset. It is well understood that in this overparameterized regime, DNNs are highly expressive and have the capacity to (over)fit arbitrary training datasets including pure noise [56]. Mysteriously however neural network models trained via simple algorithms such as stochastic gradient descent continue to predict well on yet unseen test data. In such over-parametrized scenarios there maybe infinitely many globally optimal network parameters consistent with the training data, the key challenge is to understand which network parameters (stochastic) gradient descent converges to and what are its properties. Indeed, a recent series of papers [16,52,56], suggest that solutions found by first order methods tend to have favorable generalization properties. As DNNs begin to be deployed in safety critical applications, the need for foundational understanding of their noise robustness and their unique prediction capabilities intensifies.

This paper focuses on an intriguing phenomena: overparameterized neural networks are surprisingly robust to label noise when first order methods with early stopping is used to train them. To observe this phenomena consider Figure 1 where we perform experiments on the MNIST data set. Here, we corrupt a fraction of the labels of the training data by assigning their label uniformly at random. We then fit a four layer model via stochastic gradient descent and plot various performance metrics in Figures 1a and 1b. Figure 1a (blue curve)

<sup>\*</sup>Department of Computer Science and Engineering, University of California, Riverside, CA

<sup>&</sup>lt;sup>†</sup>Ming Hsieh Department of Electrical Engineering, University of Southern California, Los Angeles, CA

<sup>&</sup>lt;sup>‡</sup>Department of Electrical and Computer Engineering, University of California, Riverside, CA

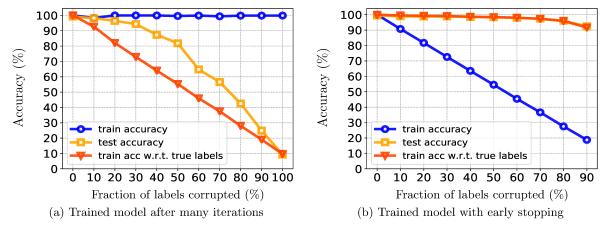


Figure 1: In these experiments we use a 4 layer neural network consisting of two convolution layers followed by two fully-connected layers to train a data set of 50,000 samples from MNIST with various amounts of random corruption on the lables. In this architecture the convolutional layers have width 64 and 128 kernels, and the fully-connected layers have 256 and 10 outputs, respectively. Overall, there are 4.8 million trainable parameters. We depict the training accuracy both w.r.t. the corrupted and uncorrupted labels as well as the test accuracy. (a) Shows the performance after 200 epochs of Adadelta where near perfect fitting to the corrupted data is achieved. (b) Shows the performance with early stopping. We observe that with early stopping the trained neural network is robust to label corruption.

shows that indeed with a sufficiently large number of iterations the neural network does in fact perfectly fit the corrupted training data. However, Figure 1a also shows that such a model does not generalize to the test data (yellow curve) and the accuracy with respect to the ground truth labels degrades (orange curve). These plots clearly demonstrate that the model overfits with many iterations. In Figure 1b we repeat the same experiment but this time stop the updates after a few iterations (i.e. use early stopping). In this case the train accuracy degrades linearly (blue curve). However, perhaps unexpected, the test accuracy (yellow curve) remains high even with a significant amount of corruption. This suggests that with early stopping the model does not overfit and generalizes to new test data. Even more surprising, the train accuracy (orange curve) with respect to the ground truth labels continues to stay around %100 even when %50 of the labels are corrupted. That is, with early stopping overparameterized neural networks even correct the corrupted labels! These plots collectively demonstrate that overparameterized neural networks when combined with early stopping have unique generalization and robustness capabilities. As we detail further in Section 4 this phenomena holds (albeit less pronounced) for reacher data models and architectures.

This paper aims to demystify the surprising robustness of overparameterized neural networks when early stopping is used. We show that gradient descent is indeed provably robust to noise/corruption on a constant fraction of the labels in such over-parametrized learning scenarios. In particular, under a fairly expressive dataset model and focusing on one-hidden layer networks, we show that after a few iterations (a.k.a. early stopping), gradient descent finds a model (i) that is within a small neighborhood of the point of initialization and (ii) only fits to the correct labels essentially ignoring the noisy labels. We complement these findings by proving that if the network is trained to overfit to the noisy labels, then the solution found by gradient descent must stray rather far from the initial model. Together, these results highlight the key features of a solution that generalizes well vs a solution that fits well.

Our theoretical results further highlight the role of the distance between final and initial network weights as a key feature that determines noise robustness vs. overfitting. This is inherently connected to the commonly used early stopping heuristic for DNN training as this heuristic helps avoid models that are too far from the point of initialization. In the presence of label noise, we show that gradient descent implicitly ignores the noisy labels as long as the model parameters remain close to the initialization. Hence, our results help explain why early stopping improves robustness and helps prevent overfitting. Under proper normalization,

the required distance between the final and initial network and the predictive accuracy of the final network is independent of the size of the network such as number of hidden nodes. Our extensive numerical experiments corroborate our theory and verify the surprising robustness of DNNs to label noise. Finally, we would like to note that while our results show that solutions found by gradient descent are inherently robust to label noise, specialized techniques such as  $\ell_1$  penalization or sample reweighting are known to further improve robustness. Our theoretical framework may enable more rigorous understandings of the benefits of such heuristics when training overparameterized models.

#### 1.2 Prior Art

Our work is connected to recent advances on theory for deep learning as well as heuristics and theory surrounding outlier robust optimization.

Robustness to label corruption: DNNs have the ability to fit to pure noise [56], however they are also empirically observed to be highly resilient to label noise and generalize well despite large corruption [44] In addition to early stopping, several heuristics have been proposed to specifically deal with label noise [26,30,36,42,47,57]. See also [23,37,43,48] for additional work on dealing with label noise in classification tasks. When learning from pairwise relations, noisy labels can be connected to graph clustering and community detection problems [1,14,54]. Label noise is also connected to outlier robustness in regression which is a traditionally well-studied topic. In the context of robust regression and high-dimensional statistics, much of the focus is on regularization techniques to automatically detect and discard outliers by using tools such as  $\ell_1$ penalization [6, 10, 15, 17, 22, 32, 35]. We would also like to note that there is an interesting line of work that focuses on developing robust algorithms for corruption not only in the labels but also input data [19, 31, 41]. Overparameterized neural networks: Intriguing properties and benefits of overparameterized neural networks has been the focus of a growing list of publications [4,11,12,18,28,49,51,53,56,58]. A recent line of work [2,3,20,21,33,38,59] show that overparameterized neural networks can fit the data with random initialization if the number of hidden nodes are polynomially large in the size of the dataset. Recently in [40] we showed that this conclusion continues to hold with more modest amounts of overparameterization and as soon as the number of parameters of the model exceed the square of the size of the training data set. This line of work however is not informative about the robustness of the trained network against corrupted labels. Indeed, such theory predicts that (stochastic) gradient descent will eventually fit the corrupted labels. In contrast, our focus here is not in finding a global minima, rather a solution that is robust to label corruption. In particular, we show that with early stopping we fit to the correct labels without overfitting to the corrupted training data. Our result also defers from this line of research in another way. The key property utilized in this research area is that the Jacobian of the neural network is well-conditioned at a random initialization if the dataset is sufficiently diverse (e.g. if the points are well-separated). In contrast, in our model the Jacobian is inherently low-rank with the rank of the Jacobian corresponding to different clusters/classes within the dataset. We harness this low-rank nature to prove that gradient descent is robust to label corruptions. We further utilize this low-rank structure to explain why neural networks can work with much more modest amounts of overparameterization where the number of parameters in the model exceeds the number of clusters raised to the fourth power and is independent of the number of data points. Furthermore, our numerical experiments verify that the Jacobian matrix of real datasets (such as CIFAR10) indeed exhibit low-rank structure. This is closely related to the observations on the Hessian of deep networks which is empirically observed to be low-rank [45]. We would also like to note that the importance of the Jacobian for overparameterized neural network analysis has also been noted by other papers including [21, 39, 49] and also [16,29] which investigate the optimization landscape and properties of SGD for training neural networks. An equally important question to understanding the convergence behavior of optimization algorithms for overparameterized models is understanding their generalization capabilities. This is the subject of a few interesting recent papers [5,7–9,13,24,34,50]. While in this paper we do not tackle generalization in the traditional sense, we do show that solution found by gradient descent are robust to label noise/corruption which demonstrates their predictive capabilities and in turn suggests better generalization.

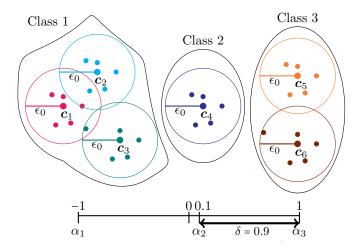


Figure 2: Visualization of the input/label samples and classes according to the clusterable dataset model in Definition 1.1. In the depicted example there are K=6 clusters,  $\bar{K}=3$  classes. In this example the number of data points is n=30 with each cluster containing 5 data points. The labels associated to classes 1, 2, and 3 are  $\alpha_1=-1$ ,  $\alpha_2=0.1$ , and  $\alpha_3=1$ , respectively so that  $\delta=0.9$ . We note that the placement of points are exaggerated for clarity. In particular, per definition the cluster center and data points all have unit Euclidean norm. Also, there is no explicit requirements that the cluster centers be separated. The depicted separation is for exposition purposes only.

#### 1.3 Models

We first describe the dataset model used in our theoretical results. In this model we assume that the input samples  $x_1, x_2, \ldots, x_n \in \mathbb{R}^d$  come from K clusters which are located on the unit Euclidian ball in  $\mathbb{R}^d$ . We also assume our data set consists of  $\overline{K} \leq K$  classes where each class can be composed of multiple clusters. We consider a deterministic data set with n samples with roughly balanced clusters each consisting on the order of n/K samples. Finally, while we allow for multiple classes, in our model we assume the labels are scalars and take values in [-1,1] interval. We formally define our dataset model below and provide an illustration in Figure 2.

**Definition 1.1 (Clusterable dataset)** Consider a data set of size n consisting of input/label pairs  $\{(\boldsymbol{x}_i,y_i)\}_{i=1}^n \in \mathbb{R}^d \times \mathbb{R}$ . We assume the input data have unit Euclidean norm and originate from K clusters with the  $\ell$ th cluster containing  $n_\ell$  data points. We assume the number of points originating from each cluster is well-balanced in the sense that  $c_{low} \frac{n}{K} \leq n_\ell \leq c_{up} \frac{n}{K}$  with  $c_{low}$  and  $c_{up}$  two numerical constants obeying  $0 < c_{low} < c_{up} < 1$ . We use  $\{\boldsymbol{c}_\ell\}_{\ell=1}^K \subset \mathbb{R}^d$  to denote the cluster centers which are distinct unit Euclidian norm vectors. We assume the input data points  $\boldsymbol{x}$  that belong to the  $\ell$ -th cluster obey

$$\|\boldsymbol{x} - \boldsymbol{c}_{\ell}\|_{\ell_2} \leq \varepsilon_0,$$

with  $\varepsilon_0 > 0$  denoting the input noise level.

We assume the labels  $y_i$  belong to one of  $\bar{K} \leq K$  classes. Specifically, we assume  $y_i \in \{\alpha_1, \alpha_2, \dots, \alpha_{\bar{K}}\}$  with  $\{\alpha_\ell\}_{\ell=1}^{\bar{K}} \in [-1, 1]$  denoting the labels associated with each class. We assume all the elements of the same cluster belong to the same class and hence have the same label. However, a class can contain multiple clusters. Finally, we assume the labels are separated in the sense that

$$|\alpha_r - \alpha_s| \ge \delta \quad \text{for} \quad r \ne s,$$
 (1.1)

with  $\delta > 0$  denoting the class separation.

<sup>&</sup>lt;sup>1</sup>This is for ease of exposition rather than a particular challenge arising in the analysis.

In the data model above  $\{c_\ell\}_{\ell=1}^K$  are the K cluster centers that govern the input distribution. We note that in this model different clusters can be assigned to the same label. Hence, this setup is rich enough to model data which is not linearly separable: e.g. over  $\mathbb{R}^2$ , we can assign cluster centers (0,1) and (0,-1) to label 1 and cluster centers (1,0) and (-1,0) to label -1. Note that the maximum number of classes are dictated by the separation  $\delta$ . In particular, we can have at most  $\overline{K} \leq \frac{2}{\delta} + 1$  classes. We remark that this model is related to the setup of [33] which focuses on providing polynomial guarantees for learning shallow networks. Finally, note that, we need some sort of separation between the cluster centers to distinguish them. While Definition 1.1 doesn't specifies such separation explicitly, Definition 2.1 establishes a notion of separation in terms of how well a neural net can distinguish the cluster centers. Next, we introduce our noisy/corrupted dataset model.

**Definition 1.2**  $((\rho, \varepsilon_0, \delta)$  **corrupted dataset**) Let  $\{(\boldsymbol{x}_i, \widetilde{y}_i)\}_{i=1}^n$  be an  $(\varepsilon_0, \delta)$  clusterable dataset with  $\alpha_1$ ,  $\alpha_2, \ldots, \alpha_{\bar{K}}$  denoting the  $\bar{K}$  possible class labels. A  $(\rho, \varepsilon_0, \delta)$  noisy/corrupted dataset  $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$  is generated from  $\{(\boldsymbol{x}_i, \widetilde{y}_i)\}_{i=1}^n$  as follows. For each cluster  $1 \leq \ell \leq K$ , at most  $s_\ell \leq \rho n_\ell$  of the labels associated with that cluster (which contains  $n_\ell$  points) is assigned to another label value chosen from  $\{\alpha_\ell\}_{\ell=1}^{\bar{K}}$ . We shall refer to the initial labels  $\{\widetilde{y}_i\}_{i=1}^n$  as the ground truth labels.

We note that this definition allows for a fraction  $\rho$  of corruptions in each cluster.

Network model: We will study the ability of neural networks to learn this corrupted dataset model. To proceed, let us introduce our neural network model. We consider a network with one hidden layer that maps  $\mathbb{R}^d$  to  $\mathbb{R}$ . Denoting the number of hidden nodes by k, this network is characterized by an activation function  $\phi$ , input weight matrix  $\mathbf{W} \in \mathbb{R}^{k \times d}$  and output weight vector  $\mathbf{v} \in \mathbb{R}^k$ . In this work, we will fix output  $\mathbf{v}$  to be a unit vector where half the entries are  $1/\sqrt{k}$  and other half are  $-1/\sqrt{k}$  to simplify exposition.<sup>2</sup> We will only optimize over the weight matrix  $\mathbf{W}$  which contains most of the network parameters and will be shown to be sufficient for robust learning. We will also assume  $\phi$  has bounded first and second order derivatives, i.e.  $|\phi'(z)|, |\phi''(z)| \le \Gamma$  for all z. The network's prediction at an input sample  $\mathbf{x}$  is given by

$$x \mapsto f(\mathbf{W}, \mathbf{x}) = \mathbf{v}^T \phi(\mathbf{W}\mathbf{x}),$$
 (1.2)

where the activation function  $\phi$  applies entrywise. Given a dataset  $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$ , we shall train the network via minimizing the empirical risk over the training data via a quadratic loss

$$\mathcal{L}(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^{n} (y_i - f(\mathbf{x}_i, \mathbf{W}))^2.$$
 (1.3)

In particular, we will run gradient descent with a constant learning rate  $\eta$ , starting from a random initialization  $W_0$  via the following updates

$$\mathbf{W}_{\tau+1} = \mathbf{W}_{\tau} - \eta \nabla \mathcal{L}(\mathbf{W}_{\tau}). \tag{1.4}$$

## 2 Main results

#### 2.1 Robustness of neural network to label noise with early stopping

Our main result shows that overparameterized neural networks, when trained via gradient descent using early stopping are fairly robust to label noise. The ability of neural networks to learn from the training data, even without label corruption, naturally depends on the diversity of the input training data. Indeed, if two input data are nearly the same but have different uncorrupted labels reliable learning is difficult. We will quantify this notion of diversity via a notion of condition number related to a covariance matrix involving the activation  $\phi$  and the cluster centers  $\{c_\ell\}_{\ell=1}^K$ .

 $<sup>^{2}</sup>$ If the number of hidden units is odd we set one entry of  $\boldsymbol{v}$  to zero.

**Definition 2.1 (Neural Net Cluster Covariance and Condition Number)** Define the matrix of cluster centers

$$C = [c_1 \dots c_K]^T \in \mathbb{R}^{K \times d}$$
.

Let  $\mathbf{g} \sim \mathcal{N}(0, \mathbf{I}_d)$ . Define the neural net covariance matrix  $\mathbf{\Sigma}(\mathbf{C})$  as

$$\Sigma(C) = (CC^T) \bigcirc \mathbb{E}_q[\phi'(Cg)\phi'(Cg)^T].$$

Here  $\odot$  denotes the elementwise product. Also denote the minimum eigenvalue of  $\Sigma(C)$  by  $\lambda(C)$  and define the following condition number associated with the cluster centers C

$$\kappa(C) = \sqrt{\frac{d}{K}} \frac{\|C\|}{\lambda(C)}.$$

One can view  $\Sigma(C)$  as an empirical kernel matrix associated with the network where the kernel is given by  $K(c_i, c_j) = \Sigma_{ij}(C)$ . Note that  $\Sigma(C)$  is trivially rank deficient if there are two cluster centers that are identical. In this sense, the minimum eigenvalue of  $\Sigma(C)$  will quantify the ability of the neural network to distinguish between distinct cluster centers. Therefore, one can think of  $\kappa(C)$  as a condition number associated with the neural network which characterizes the distinctness/diversity of the cluster centers. The more distinct the cluster centers, the larger  $\lambda(C)$  and smaller the condition number  $\kappa(C)$  is. Indeed, based on results in [40] when the cluster centers are maximally diverse e.g. uniformly at random from the unit sphere  $\kappa(C)$  scales like a constant. Throughout we shall assume that  $\lambda(C)$  is strictly positive (and hence  $\kappa(C) < \infty$ ). This property is empirically verified to hold in earlier works [55] when  $\phi$  is a standard activation (e.g. ReLU, softplus). As a concrete example, for ReLU activation, using results from [40] one can show if the cluster centers are separated by a distance  $\nu > 0$ , then  $\lambda(C) \ge \frac{\nu}{100K^2}$ . We note that variations of the  $\lambda(C) > 0$  assumption based on the data points (i.e.  $\lambda(X) > 0$  not cluster centers) [20, 21, 40] are utilized to provide convergence guarantees for DNNs. Also see [3, 59] for other publications using related definitions.

Now that we have a quantitative characterization of distinctiveness/diversity in place we are now ready to state our main result. Throughout we use  $c_{\Gamma}$ ,  $C_{\Gamma}$ , etc. to denote constants only depending on  $\Gamma$ . We note that this Theorem is slightly simplified by ignoring logarithmic terms and precise dependencies on  $\Gamma$ . We refer the reader to Theorem 6.13 for precise statement including logarithmic terms.

Theorem 2.2 (Robust learning with early stopping-simplified) Consider an  $(s, \varepsilon_0, \delta)$  clusterable corrupted data set of input/label pairs  $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n \in \mathbb{R}^d \times \mathbb{R}$  per Definition 1.2 with cluster centers  $\{\boldsymbol{c}_\ell\}_{\ell=1}^K$  aggregated as rows of a matrix  $\boldsymbol{C} \in \mathbb{R}^{K \times d}$ . Furthermore, let  $\{\widetilde{y}_i\}_{i=1}^n$  be the corresponding uncorrupted ground truth labels. Also consider a one-hidden layer neural network with k hidden units and one output of the form  $\boldsymbol{x} \mapsto \boldsymbol{v}^T \phi(\boldsymbol{W} \boldsymbol{x})$  with  $\boldsymbol{W} \in \mathbb{R}^{k \times d}$  and  $\boldsymbol{v} \in \mathbb{R}^k$  the input-to-hidden and hidden-to-output weights. Also suppose the activation  $\phi$  obeys  $|\phi(0)| \leq \Gamma$  and  $|\phi'(z)|, |\phi''(z)| \leq \Gamma$  for all z and some  $\Gamma \geq 1$ . Furthermore, we set half of the entries of  $\boldsymbol{v}$  to  $1/\sqrt{k}$  and the other half to  $-1/\sqrt{k^3}$  and train only over  $\boldsymbol{W}$ . Starting from an initial weight matrix  $\boldsymbol{W}_0$  selected at random with i.i.d.  $\mathcal{N}(0,1)$  entries we run Gradient Descent (GD) updates of the form  $\boldsymbol{W}_{\tau+1} = \boldsymbol{W}_{\tau} - \eta \nabla \mathcal{L}(\boldsymbol{W}_{\tau})$  on the least-squares loss (1.3) with step size  $\eta = \bar{c}_{\Gamma} \frac{K}{n} \frac{1}{\|\boldsymbol{C}\|^2}$  with  $\bar{c}_{\Gamma}$ . Furthermore, assume the number of parameters obey

$$kd \geq C_{\Gamma} \kappa^4(\mathbf{C}) \frac{K^4}{d},$$

with  $\kappa(C)$  the neural net cluster condition number pre Definition 2.1. Then as long as  $\epsilon_0 \leq \widetilde{c}_{\Gamma}/K^2$  and  $\rho \leq \frac{\delta}{8}$  with probability at least  $1 - 3/K^{100}$ , after  $\tau_0 = c_{\Gamma} \frac{K}{d} \lambda(C) \kappa^2(C) \log(\frac{1}{\rho})$  iterations, the neural network  $f(\cdot, \mathbf{W}_{\tau_0})$  found by gradient descent assigns all the input samples  $\mathbf{x}_i$  to the correct ground truth labels  $\widetilde{y}_i$ . That is,

$$\arg\min_{\alpha_{\ell}:1\leq \ell\leq \bar{K}}|f(\boldsymbol{W}_{\tau},\boldsymbol{x}_{i})-\alpha_{\ell}|=\widetilde{y}_{i},\tag{2.1}$$

<sup>&</sup>lt;sup>3</sup>If k is odd we set one entry to zero  $\lfloor \frac{k-1}{2} \rfloor$  to  $1/\sqrt{k}$  and  $\lfloor \frac{k-1}{2} \rfloor$  entries to  $-1/\sqrt{k}$ .

holds for all  $1 \le i \le n$ . Furthermore, for all  $0 \le \tau \le \tau_0$ , the distance to the initial point obeys

$$\|\boldsymbol{W}_{\tau} - \boldsymbol{W}_{0}\|_{F} \leq \bar{C}_{\Gamma} \left( \sqrt{K} + \frac{K^{2}}{\|\boldsymbol{C}\|^{2}} \tau \varepsilon_{0} \right).$$

Theorem 2.2 shows that gradient descent with early stopping has a few intriguing properties. We further discuss these properties below.

**Robustness.** The solution found by gradient descent with early stopping degrades gracefully as the label corruption level  $\rho$  grows. In particular, as long as  $\rho \leq \delta/8$ , the final model is able to correctly classify all samples including the corrupted ones. In our setup, intuitively label gap obeys  $\delta \sim \frac{1}{K}$ , hence, we prove robustness to

Total Number of corrupted labels 
$$\lesssim \frac{n}{\bar{k}}$$
.

This result is independent of number of clusters and only depends on number of classes. An interesting future direction is to improve this result to allow on the order of n corrupted labels. Such a result maybe possible by using a multi-output classification neural network.

Early stopping time. We show that gradient descent finds a model that is robust to outliers after a few iterations. In particular using the maximum allowed step size, the required number of iterations is of the order of  $\frac{K}{d}\lambda(C)\kappa^2(C)\log(\frac{1}{\rho})$  which scales with K/d up to condition numbers.

Modest overparameterization. Our result requires modest overparemetrization and apply as soon as the

Modest overparameterization. Our result requires modest overparemetrization and apply as soon as the number of parameters exceed the number of classes to the power four  $(kd \gtrsim K^4)$ . Interestingly, under our data model the required amount of overparameterization is essentially independent of the size of the training data n(ignoring logarithmic terms) and conditioning of the data points, only depending on the number of clusters and conditioning of the cluster centers. This can be interpreted as ensuring that the network has enough capacity to fit the cluster centers  $\{c_\ell\}_{\ell=1}^K$  and the associated true labels.

**Distance from initialization.** Another feature of Theorem 2.2 is that the network weights do not stray far from the initialization as the distance between the initial model and the final model (at most) grows with the square root of the number of clusters  $(\sqrt{K})$ . This  $\sqrt{K}$  dependence implies that the more clusters there are, the updates travel further away but continue to stay within a certain radius. This dependence is intuitive as the *Rademacher complexity* of the function space is dictated by the distance to initialization and should grow with the square-root of the number of input clusters to ensure the model is expressive enough to learn the dataset.

Before we end this section we would like to note that in the limit of  $\epsilon_0 \to 0$  where the input data set is perfectly clustered one can improve the amount of overparamterization. Indeed, the result above is obtained via a perturbation argument from this more refined result stated below.

Theorem 2.3 (Training with perfectly clustered data) Consier the setting and assumptions of Theorem 2.3 with  $\epsilon_0 = 0$ . Starting from an initial weight matrix  $\mathbf{W}_0$  selected at random with i.i.d.  $\mathcal{N}(0,1)$  entries we run Gradient Descent (GD) updates of the form  $\mathbf{W}_{\tau+1} = \mathbf{W}_{\tau} - \eta \nabla \mathcal{L}(\mathbf{W}_{\tau})$  on the least-squares loss (1.3) with step size  $\eta \leq \frac{K}{2c_{up}n\Gamma^2\|\mathbf{C}\|^2}$ . Furthermore, assume the number of parameters obey

$$kd \ge C\Gamma^4 \kappa^2(\boldsymbol{C})K^2,$$

with  $\kappa(C)$  the neural net cluster condition number per Definition 2.1. Then, with probability at least  $1-2/K^{100}$  over randomly initialized  $W_0 \overset{i.i.d.}{\sim} \mathcal{N}(0,1)$ , the iterates  $W_{\tau}$  obey the following properties.

ullet The distance to initial point  $oldsymbol{W}_0$  is upper bounded by

$$\|\boldsymbol{W}_{\tau} - \boldsymbol{W}_{0}\|_{F} \le c\Gamma \sqrt{\frac{K \log K}{\lambda(\boldsymbol{C})}}.$$

• After  $\tau \geq \tau_0 := c \frac{K}{\eta n \lambda(C)} \log \left( \frac{\Gamma \sqrt{n \log K}}{\rho} \right)$  iterations, the entrywise predictions of the learned network with respect to the ground truth labels  $\{\widetilde{y}_i\}_{i=1}^n$  satisfy

$$|f(\mathbf{W}_{\tau}, \mathbf{x}_i) - \widetilde{y}_i| \leq 4\rho$$

for all  $1 \le i \le n$ . Furthermore, if the noise level  $\rho$  obeys  $\rho \le \delta/8$  the network predicts the correct label for all samples i.e.

$$\arg\min_{\alpha_{\ell}:1\leq\ell\leq\bar{K}}|f(\boldsymbol{W}_{\tau},\boldsymbol{x}_{i})-\alpha_{\ell}|=\widetilde{y}_{i}\quad for\quad i=1,2,\ldots,n.$$
(2.2)

This result shows that in the limit  $\epsilon_0 \to 0$  where the data points are perfectly clustered, the required amount of overparameterization can be reduced from  $kd \gtrsim K^4$  to  $kd \gtrsim K^2$ . In this sense this can be thought of a nontrivial analogue of [40] where the number of data points are replaced with the number of clusters and the condition number of the data points is replaced with a cluster condition number. This can be interpreted as ensuring that the network has enough capacity to fit the cluster centers  $\{c_\ell\}_{\ell=1}^K$  and the associated true labels. Interestingly, the robustness benefits continue to hold in this case. However, in this perfectly clustered scenario there is no need for early stopping and a robust network is trained as soon as the number of iterations are sufficiently large. Infact, in this case given the clustered nature of the input data the network never overfits to the corrupted data even after many iterations.

## 2.2 To (over)fit to corrupted labels requires straying far from initialization

In this section we wish to provide further insight into why early stopping enables robustness and generalizable solutions. Our main insight is that while a neural network maybe expressive enough to fit a corrupted dataset, the model has to travel a longer distance from the point of initialization as a function of the distance from the cluster centers  $\varepsilon_0$  and the amount of corruption. We formalize this idea as follows. Suppose

- 1. two input points are close to each other (e.g. they are from the same cluster),
- 2. but their labels are different, hence the network has to map them to distant outputs.

Then, the network has to be large enough so that it can amplify the small input difference to create a large output difference. Our first result formalizes this for a randomly initialized network. Our random initialization picks W with i.i.d. standard normal entries which ensures that the network is isometric i.e. given input x,  $\mathbb{E}[f(W,x)^2] = \mathcal{O}(\|x\|_{\ell_2}^2)$ .

**Theorem 2.4** Let  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$  be two vectors with unit Euclidean norm obeying  $\|\mathbf{x}_2 - \mathbf{x}_1\|_{\ell_2} \leq \epsilon_0$ . Let  $f(\mathbf{W}, \mathbf{x}) = \mathbf{v}^T \phi(\mathbf{W}\mathbf{x})$  where  $\mathbf{v}$  is fixed,  $\mathbf{W} \in \mathbb{R}^{k \times d}$ , and  $k \geq cd$  with c > 0 a fixed constant. Assume  $|\phi'|, |\phi''| \leq \Gamma$ . Let  $y_1$  and  $y_2$  be two scalars satisfying  $|y_2 - y_1| \geq \delta$ . Suppose  $\mathbf{W}_0 \stackrel{i.i.d.}{\sim} \mathcal{N}(0,1)$ . Then, with probability at least  $1 - 2e^{-(k+d)} - 2e^{-\frac{t^2}{2}}$ , for any  $\mathbf{W} \in \mathbb{R}^{k \times d}$  such that  $\|\mathbf{W} - \mathbf{W}_0\|_F \leq c\sqrt{k}$  and

$$f(\boldsymbol{W}, \boldsymbol{x}_1) = y_1$$
 and  $f(\boldsymbol{W}, \boldsymbol{x}_2) = y_2$ ,

holds, we have

$$\|\boldsymbol{W} - \boldsymbol{W}_0\| \ge \frac{\delta}{C\Gamma\varepsilon_0} - \frac{t}{1000}.$$

In words, this result shows that in order to fit to a data set with a *single corrupted label*, a randomly initialized network has to traverse a distance of at least  $\delta/\varepsilon_0$ . The next lemma clarifies the role of the corruption amount s and shows that more label corruption within a fixed class requires a model with a larger norm in order to fit the labels. For this result we consider a randomized model with  $\varepsilon_0^2$  input noise variance.

**Lemma 2.5** Let  $c \in \mathbb{R}^d$  be a cluster center. Consider 2s data points  $\{x_i\}_{i=1}^s$  and  $\{\widetilde{x}_i\}_{i=1}^s$  in  $\mathbb{R}^d$  generated i.i.d. around c according to the following distribution

$$c + g$$
 with  $g \sim \mathcal{N}(0, \frac{\varepsilon_0^2}{d} I_d)$ .

Assign  $\{\boldsymbol{x}_i\}_{i=1}^s$  with labels  $y_i = y$  and  $\{\widetilde{\boldsymbol{x}}_i\}_{i=1}^s$  with labels  $\widetilde{y}_i = \widetilde{y}$  and assume these two labels are  $\delta$  separated i.e.  $|y - \widetilde{y}| \ge \delta$ . Also suppose  $s \le d$  and  $|\phi'| \le \Gamma$ . Then, any  $\boldsymbol{W} \in \mathbb{R}^{k \times d}$  satisfying

$$f(\mathbf{W}, \mathbf{x}_i) = y_i$$
 and  $f(\mathbf{W}, \widetilde{\mathbf{x}}_i) = \widetilde{y}_i$  for  $i = 1, \dots, s$ ,

obeys  $\|\mathbf{W}\|_F \ge \frac{\sqrt{s}\delta}{5\Gamma\varepsilon_0}$  with probability at least  $1 - e^{-d/2}$ .

Unlike Theorem 2.4 this result lower bounds the network norm in lieu of the distance to the initialization  $W_0$ . However, using the triangular inequality we can in turn get a guarantee on the distance from initialization  $W_0$  via triangle inequality as long as  $\|W_0\|_F \lesssim \mathcal{O}(\sqrt{s\delta/\varepsilon_0})$  (e.g. by choosing a small  $\varepsilon_0$ ).

The above Theorem implies that the model has to traverse a distance of at least

$$\|\boldsymbol{W}_{\tau} - \boldsymbol{W}_{0}\|_{F} \gtrsim \sqrt{\frac{\rho n}{K}} \frac{\delta}{\varepsilon_{0}},$$

to perfectly fit corrupted labels. In contrast, we note that the conclusions of the upper bound in Theorem 2.2 show that to be able to fit to the uncorrupted true labels the distance to initialization grows at most by  $\tau \varepsilon_0$  after  $\tau$  iterates. This demonstrates that there is a gap in the required distance to initialization for fitting enough to generalize and overfitting. To sum up, our results highlight that, one can find a network with good generalization capabilities and robustness to label corruption within a small neighborhood of the initialization and that the size of this neighborhood is independent of the corruption. However, to fit to the corrupted labels, one has to travel much more, increasing the search space and likely decreasing generalization ability. Thus, early stopping can enable robustness without overfitting by restricting the distance to the initialization.

## 3 Technical Approach and General Theory

In this section, we outline our approach to proving robustness of overparameterized neural networks. Towards this goal, we consider a general formulation where we aim to fit a general nonlinear model of the form  $x \mapsto f(\theta, x)$  with  $\theta \in \mathbb{R}^p$  denoting the parameters of the model. For instance in the case of neural networks  $\theta$  represents its weights. Given a data set of n input/label pairs  $\{(x_i, y_i)\}_{i=1}^n \subset \mathbb{R}^d \times \mathbb{R}$ , we fit to this data by minimizing a nonlinear least-squares loss of the form

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^{n} (y_i - f(\boldsymbol{\theta}, \boldsymbol{x}_i))^2.$$

which can also be written in the more compact form

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \| f(\boldsymbol{\theta}) - \boldsymbol{y} \|_{\ell_2}^2 \quad \text{with} \quad f(\boldsymbol{\theta}) \coloneqq \begin{bmatrix} f(\boldsymbol{\theta}, \boldsymbol{x}_1) \\ f(\boldsymbol{\theta}, \boldsymbol{x}_2) \\ \vdots \\ f(\boldsymbol{\theta}, \boldsymbol{x}_n) \end{bmatrix}.$$

To solve this problem we run gradient descent iterations with a constant learning rate  $\eta$  starting from an initial point  $\theta_0$ . These iterations take the form

$$\theta_{\tau+1} = \theta_{\tau} - \eta \nabla \mathcal{L}(\theta_{\tau}) \quad \text{with} \quad \nabla \mathcal{L}(\theta) = \mathcal{J}^{T}(\theta) (f(\theta) - y).$$
 (3.1)

Here,  $\mathcal{J}(\boldsymbol{\theta})$  is the  $n \times p$  Jacobian matrix associated with the nonlinear mapping f defined via

$$\mathcal{J}(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial f(\boldsymbol{\theta}, \boldsymbol{x}_1)}{\partial \boldsymbol{\theta}} & \dots & \frac{\partial f(\boldsymbol{\theta}, \boldsymbol{x}_n)}{\partial \boldsymbol{\theta}} \end{bmatrix}^T.$$
(3.2)

#### 3.1 Bimodal jacobian structure

Our approach is based on the hypothesis that the nonlinear model has a Jacobian matrix with bimodal spectrum where few singular values are large and remaining singular values are small. This assumption is inspired by the fact that realistic datasets are clusterable in a proper, possibly nonlinear, representation space. Indeed, one may argue that one reason for using neural networks is to automate the learning of such a representation (essentially the input to the softmax layer). We formalize the notion of bimodal spectrum below.

**Assumption 1 (Bimodal Jacobian)** Let  $\beta \geq \alpha \geq \epsilon > 0$  be scalars. Let  $f : \mathbb{R}^p \to \mathbb{R}^n$  be a nonlinear mapping and consider a set  $\mathcal{D} \subset \mathbb{R}^p$  containing the initial point  $\theta_0$  (i.e.  $\theta_0 \in \mathcal{D}$ ). Let  $\mathcal{S}_+ \subset \mathbb{R}^n$  be a subspace and  $\mathcal{S}_-$  be its complement. We say the mapping f has a Bimodal Jacobian with respect to the complementary subpspaces  $\mathcal{S}_+$  and  $\mathcal{S}_-$  as long as the following two assumptions hold for all  $\theta \in \mathcal{D}$ .

• Spectrum over  $S_+$ : For all  $v \in S_+$  with unit Euclidian norm we have

$$\alpha \le \|\mathcal{J}^T(\boldsymbol{\theta})\boldsymbol{v}\|_{\ell_2} \le \beta.$$

• Spectrum over  $S_{-}$ : For all  $v \in S_{-}$  with unit Euclidian norm we have

$$\|\mathcal{J}^T(\boldsymbol{\theta})\boldsymbol{v}\|_{\ell_2} \leq \epsilon.$$

We will refer to  $S_+$  as the signal subspace and  $S_-$  as the noise subspace.

When  $\epsilon \ll \alpha$  the Jacobian is approximately low-rank. An extreme special case of this assumption is where  $\epsilon = 0$  so that the Jacobian matrix is exactly low-rank. We formalize this assumption below for later reference.

**Assumption 2 (Low-rank Jacobian)** Let  $\beta \geq \alpha > 0$  be scalars. Consider a set  $\mathcal{D} \subset \mathbb{R}^p$  containing the initial point  $\theta_0$  (i.e.  $\theta_0 \in \mathcal{D}$ ). Let  $\mathcal{S}_+ \subset \mathbb{R}^n$  be a subspace and  $\mathcal{S}_-$  be its complement. For all  $\theta \in \mathcal{D}$ ,  $v \in \mathcal{S}_+$  and  $w \in \mathcal{S}_-$  with unit Euclidian norm, we have that

$$\alpha \leq \left\| \mathcal{J}^T(\boldsymbol{\theta}) \boldsymbol{v} \right\|_{\ell_2} \leq \beta \quad and \quad \left\| \mathcal{J}^T(\boldsymbol{\theta}) \boldsymbol{w} \right\|_{\ell_2} = 0.$$

Our dataset model in Definition 1.2 naturally has a low-rank Jacobian when  $\epsilon_0 = 0$  and each input example is equal to one of the K cluster centers  $\{c_\ell\}_{\ell=1}^K$ . In this case, the Jacobian will be at most rank K since each row will be in the span of  $\{\frac{\partial f(c_\ell, \theta)}{\partial \theta}\}_{\ell=1}^K$ . The subspace  $\mathcal{S}_+$  is dictated by the membership of each cluster as follows: Let  $\Lambda_\ell \subset \{1, \ldots, n\}$  be the set of coordinates i such that  $x_i = c_\ell$ . Then, subspace is characterized by

$$S_+ = \{ \boldsymbol{v} \in \mathbb{R}^n \mid \boldsymbol{v}_{i_1} = \boldsymbol{v}_{i_2} \text{ for all } i_1, i_2 \in \Lambda_\ell \text{ and } 1 \le \ell \le K \}.$$

When  $\epsilon_0 > 0$  and the data points of each cluster are not the same as the cluster center we have the bimodal Jacobian structure of Assumption 1 where over  $S_-$  the spectral norm is small but nonzero.

In Section 4, we verify that the Jacobian matrix of real datasets indeed have a bimodal structure i.e. there are few large singular values and the remaining singular values are small which further motivate Assumption 2. This is inline with earlier papers which observed that Hessian matrices of deep networks have bimodal spectrum (approximately low-rank) [45] and is related to various results demonstrating that there are flat directions in the loss landscape [27].

## 3.2 Meta result on learning with label corruption

Define the *n*-dimensional residual vector  $\mathbf{r}$  where  $\mathbf{r}(\boldsymbol{\theta}) = [f(\mathbf{x}_1, \boldsymbol{\theta}) - \mathbf{y}_1 \dots f(\mathbf{x}_n, \boldsymbol{\theta}) - \mathbf{y}_n]^T$ . A key idea in our approach is that we argue that (1) in the absence of any corruption  $\mathbf{r}(\boldsymbol{\theta})$  approximately lies on the subspace  $\mathcal{S}_+$  and (2) if the labels are corrupted by a vector  $\mathbf{e}$ , then  $\mathbf{e}$  approximately lies on the complement space. Before we state our general result we need to discuss another assumption and definition.

**Assumption 3 (Smoothness)** The Jacobian mapping  $\mathcal{J}(\boldsymbol{\theta})$  associated to a nonlinear mapping  $f: \mathbb{R}^p \to \mathbb{R}^n$  is L-smooth if for all  $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \mathbb{R}^p$  we have  $\|\mathcal{J}(\boldsymbol{\theta}_2) - \mathcal{J}(\boldsymbol{\theta}_1)\| \le L \|\boldsymbol{\theta}_2 - \boldsymbol{\theta}_1\|_{\ell_0}$ .

Additionally, to connect our results to the number of corrupted labels, we introduce the notion of subspace diffusedness defined below.

**Definition 3.1** (Diffusedness)  $S_+$  is  $\gamma$  diffused if for any vector  $\mathbf{v} \in S_+$ 

$$\|\boldsymbol{v}\|_{\ell_{\infty}} \leq \sqrt{\gamma/n} \|\boldsymbol{v}\|_{\ell_{2}},$$

holds for some  $\gamma > 0$ .

The following theorem is our meta result on the robustness of gradient descent to sparse corruptions on the labels when the Jacobian mapping is exactly low-rank. Theorem 2.3 for the perfectly clustered data ( $\epsilon_0 = 0$ ) is obtained by combining this result with specific estimates developed for neural networks.

Theorem 3.2 (Gradient descent with label corruption) Consider a nonlinear least squares problem of the form  $\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2} \| f(\boldsymbol{\theta}) - \boldsymbol{y} \|_{\ell_2}^2$  with the nonlinear mapping  $f : \mathbb{R}^p \to \mathbb{R}^n$  obeying assumptions 2 and 3 over a unit Euclidian ball of radius  $\frac{4\|\boldsymbol{r}_0\|_{\ell_2}}{\alpha}$  around an initial point  $\boldsymbol{\theta}_0$  and  $\boldsymbol{y} = [y_1 \dots y_n] \in \mathbb{R}^n$  denoting the corrupted labels. Also let  $\widetilde{\boldsymbol{y}} = [\widetilde{y}_1 \dots \widetilde{y}_n] \in \mathbb{R}^n$  denote the uncorrupted labels and  $\boldsymbol{e} = \boldsymbol{y} - \widetilde{\boldsymbol{y}}$  the corruption. Furthermore, suppose the initial residual  $f(\boldsymbol{\theta}_0) - \widetilde{\boldsymbol{y}}$  with respect to the uncorrupted labels obey  $f(\boldsymbol{\theta}_0) - \widetilde{\boldsymbol{y}} \in \mathcal{S}_+$ . Then, running gradient descent updates of the from (3.1) with a learning rate  $\eta \leq \frac{1}{2\beta^2} \min\left(1, \frac{\alpha\beta}{L\|\boldsymbol{r}_0\|_{\ell_2}}\right)$ , all iterates obey

$$\|\boldsymbol{\theta}_{\tau} - \boldsymbol{\theta}_0\|_{\ell_2} \leq \frac{4\|\boldsymbol{r}_0\|_{\ell_2}}{\alpha}.$$

Furthermore, assume  $\nu > 0$  is a precision level obeying  $\nu \geq \|\Pi_{\mathcal{S}_+}(e)\|_{\ell_\infty}$ . Then, after  $\tau \geq \frac{5}{\eta\alpha^2}\log\left(\frac{\|r_0\|_{\ell_2}}{\nu}\right)$  iterations,  $\theta_{\tau}$  achieves the following error bound with respect to the true labels

$$||f(\boldsymbol{\theta}_{\tau}) - \widetilde{\boldsymbol{y}}||_{\ell_{\infty}} \leq 2\nu.$$

Furthermore, if e has at most s nonzeros and  $S_+$  is  $\gamma$  diffused per Definition 3.1, then using  $\nu = \|\Pi_{S_+}(e)\|_{\ell_\infty}$ 

$$||f(\boldsymbol{\theta}_{\tau}) - \widetilde{\boldsymbol{y}}||_{\ell_{\infty}} \le 2||\Pi_{\mathcal{S}_{+}}(\boldsymbol{e})||_{\ell_{\infty}} \le \frac{\gamma\sqrt{s}}{n}||\boldsymbol{e}||_{\ell_{2}}.$$

This result shows that when the Jacobian of the nonlinear mapping is low-rank, gradient descent enjoys two intriguing properties. First, gradient descent iterations remain rather close to the initial point. Second, the estimated labels of the algorithm enjoy sample-wise robustness guarantees in the sense that the noise in the estimated labels are gracefully distributed over the dataset and the effects on individual label estimates are negligible. This theorem is the key result that allows us to prove Theorem 2.3 when the data points are perfectly clustered ( $\epsilon_0 = 0$ ). Furthermore, this theorem when combined with a perturbation analysis allows us to deal with data that is not perfectly clustered ( $\epsilon_0 > 0$ ) and to conclude that with early stopping neural networks are rather robust to label corruption (Theorem 2.2).

Finally, we note that a few recent publication [3,21,39] require the Jacobian to be well-conditioned to fit labels perfectly. In contrast, our low-rank model cannot perfectly fit the corrupted labels. Furthermore, when the Jacobian is bimodal (as seems to be the case for many practical data sets and neural network models) it would take a very long time to perfectly fit the labels and as demonstrated earlier such a model does not generalize and is not robust to corruptions. Instead we focus on proving robustness with early stopping.

<sup>&</sup>lt;sup>4</sup>Note that, if  $\frac{\partial \mathcal{J}(\theta)}{\partial \theta}$  is continuous, the smoothness condition holds over any compact domain (albeit for a possibly large L).

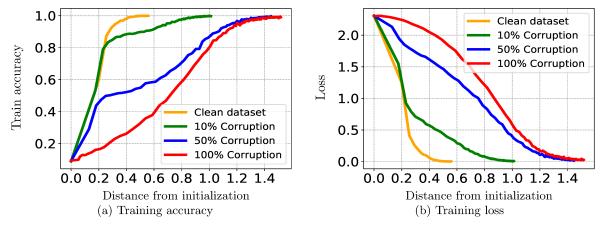


Figure 3: We depict the training accuracy of a LENET model trainined on 3000 samples from MNIST as a function of relative distance from initialization. Here, the x-axis keeps track of the distance between the current and initial weights of all layers combined.

#### 3.3 To (over)fit to corrupted labels requires straying far from initialization

In this section we state a result that provides further justification as to why early stopping of gradient descent leads to more robust models without overfitting to corrupted labels. This is based on the observation that while finding an estimate that fits the uncorrupted labels one does not have to move far from the initial estimate in the presence of corruption one has to stray rather far from the initialization with the distance from initialization increasing further in the presence of more corruption. We make this observation rigorous below by showing that it is more difficult to fit to the portion of the residual that lies on the noise space compared to the portion on the signal space (assuming  $\alpha \gg \epsilon$ ).

**Theorem 3.3** Denote the residual at initialization  $\theta_0$  by  $\mathbf{r}_0 = f(\theta_0) - \mathbf{y}$ . Define the residual projection over the signal and noise space as

$$E_{+} = \|\Pi_{S_{+}}(\mathbf{r}_{0})\|_{\ell_{2}} \quad and \quad E_{-} = \|\Pi_{S_{-}}(\mathbf{r}_{0})\|_{\ell_{2}}.$$

Suppose Assumption 1 holds over an Euclidian ball  $\mathcal{D}$  of radius  $R < \max\left(\frac{E_+}{\beta}, \frac{E_-}{\varepsilon}\right)$  around the initial point  $\theta_0$  with  $\alpha \ge \epsilon$ . Then, over  $\mathcal{D}$  there exists no  $\theta$  that achieves zero training loss. In particular, if  $\mathcal{D} = \mathbb{R}^p$ , any parameter  $\theta$  achieving zero training loss  $(f(\theta) = y)$  satisfies the distance bound

$$\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_{\ell_2} \ge \max\left(\frac{E_+}{\beta}, \frac{E_-}{\varepsilon}\right).$$

This theorem shows that the higher the corruption (and hence  $E_{-}$ ) the further the iterates need to stray from the initial model to fit the corrupted data.

## 4 Numerical experiments

We conduct several experiments to investigate the robustness capabilities of deep networks to label corruption. In our first set of experiments, we explore the relationship between loss, accuracy, and amount of label corruption on the MNIST dataset to corroborate our theory. Our next experiments study the distribution of the loss and the Jacobian on the CIFAR-10 dataset. Finally, we simulate our theoretical model by generating data according to the corrupted data model of Definition 1.2 and verify the robustness capability of gradient descent with early stopping in this model.

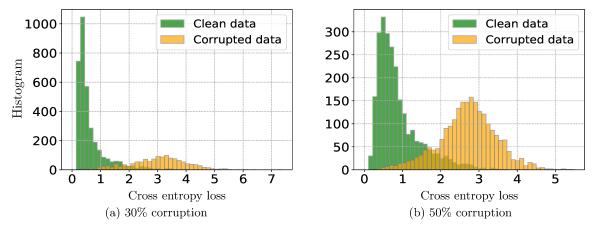


Figure 4: Histogram of the cross entropy loss of individual data points based on a model trained on 50,000 samples from CIFAR-10 with early stopping. Plot depicts 5000 random samples from these 50,000 samples. The loss distribution of clean and corrupted data are separated but gracefully overlap as the corruption level increases.

In Figure 3, we train the same model used in Figure 1 with n=3,000 MNIST samples for different amounts of corruption. Our theory predicts that more label corruption leads to a larger distance to initialization. To probe this hypothesis, Figure 3a and 3b visualizes training accuracy and training loss as a function of the distance from the initialization. These results demonstrate that the distance from initialization gracefully increase with more corruption.

Next, we study the distribution of the individual sample losses on the CIFAR-10 dataset. We conducted two experiments using Resnet-20 with cross entropy loss<sup>5</sup>. In Figure 4 we assess the noise robustness of gradient descent where we used all 50,000 samples with either 30% random corruption or 50% random corruption. Theorem 2.3 predicts that when the corruption level is small, the loss distribution of corrupted vs clean samples should be separable. Figure 4 shows that when 30% of the data is corrupted the distributions are approximately separable. When we increase the shuffling amount to 50% the training loss on the clean data increases as predicted by our theory and the distributions start to gracefully overlap.

As described in Section 3, our technical framework utilizes a bimodal prior on the Jacobian matrix (3.2) of the model. We now further investigate this hypothesis. For a multiclass task, the Jacobian matrix is essentially a 3-way tensor where dimensions are sample size (n), total number of parameters in the model (p), and the number of classes  $(\bar{K})$ . The neural network model we used for CIFAR 10 has around 270,000 parameters in total. In Figure 5 we illustrate the singular value spectrum of the two multiclass Jacobian models where we form the Jacobian from all layers except the five largest (in total we use  $\bar{p} \approx 90,000$  parameters).<sup>6</sup> We train the model with all samples and focus on the spectrum before and after the training. In Figure 5a, we picked n = 1000 samples and unfolded this tensor along parameters to obtain a  $10,000 \times 90,000$  matrix which verifies our intuition on bimodality. In particular, only 10 to 20 singular values are larger than 0.1× the top one. This is consistent with earlier works that studied the Hessian spectrum. However, focusing on the Jacobian has the added advantage of requiring only first order information [25,45]. A disadvantage is that the size of Jacobian grows with number of classes. Intuitively, cross entropy loss focuses on the class associated with the label hence in Figure 5b, we only picked the partial derivative associated with the correct class so that each sample is responsible for a single (size  $\bar{p}$ ) vector. This allowed us to scale to n = 10000 samples and the corresponding spectrum is strikingly similar. Another intriguing finding is that the spectrums of before and after training are fairly close to each other highlighting that even at random initialization, spectrum is bimodal.

 $<sup>^5</sup>$ We opted for cross entropy as it is the standard classification loss however least-squares loss achieves similar accuracy.

<sup>&</sup>lt;sup>6</sup>We depict the smaller Jacobian due to the computational cost of calculating the full Jacobian.

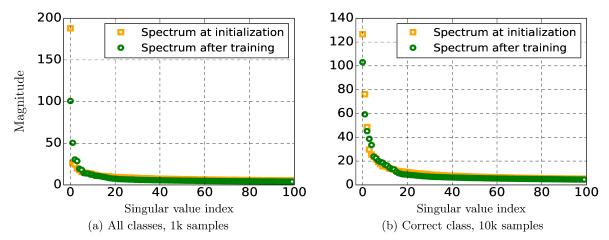


Figure 5: Spectrum of the Jacobian obtained by plotting the singular values. (a) is obtained by forming the Jacobian by taking partial derivatives of all classes associated with a sample for 1000 samples. (b) is obtained by taking the class corresponding to the label for 10000 samples.

$\# > 0.1 \times \text{top singular}$	At initialization	After training
All classes	4	14
Correct class	15	16

Table 1: Jacobian of the network has few singular values that are significantly large i.e. larger than 0.1× the spectral norm. This is true whether we consider the initial network or final network.

In Figure 6, we turn our attention to verifying our findings for the corrupted dataset model of Definition 1.2. We generated K=2 classes where the associated clusters centers are generated uniformly at random on the unit sphere of  $\mathbb{R}^{d=20}$ . We also generate the input samples at random around these two clusters uniformly at random on a sphere of radius  $\varepsilon_0=0.5$  around the corresponding cluster center. Hence, the clusters are guaranteed to be at least 1 distance from each other to prevent overlap. Overall we generate n=400 samples (200 per class/cluster). Here,  $\bar{K}=K=2$  and the class labels are 0 and 1. We picked a network with k=1000 hidden units and trained on a data set with 400 samples where 30% of the labels were corrupted. Figure 6a plots the trajectory of training error and highlights the model achieves good classification in the first few iterations and ends up overfitting later on. In Figures 6b and 6c, we focus on the loss distribution of 6a at iterations 80 and 4500. In this figure, we visualize the loss distribution of clean and corrupted data. Figure 6b highlights the loss distribution with early stopping and implies that the gap between corrupted and clean loss distributions is surprisingly resilient despite a large amount of corruption and the high-capacity of the model. In Figure 6c, we repeat plot after many more iterations at which point the model overfits. This plot shows that the distribution of the two classes overlap demonstrating that the model has overfit the corruption and lacks generalization/robustness.

## 5 Conclusions

In this paper, we studied the robustness of overparameterized neural networks to label corruption from a theoretical lens. We provided robustness guarantees for training networks with gradient descent when early stopping is used and complemented these guarantees with lower bounds. Our results point to the distance between final and initial network weights as a key feature to determine robustness vs. overfitting which is inline with weight decay and early stopping heuristics. We also carried out extensive numerical experiments to verify the theoretical predictions as well as technical assumptions. While our results shed light on the

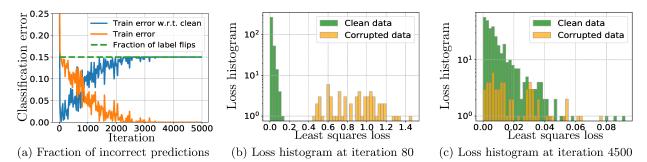


Figure 6: We experiment with the corrupted dataset model of Definition 1.2. We picked K = 2 classes and set n = 400 and  $\varepsilon_0 = 0.5$ . Trained 30% corrupted data with k = 1000 hidden units. Each corruption has 50% chance to remain in the correct class hence around 15% of the labels are actually flipped which corresponds to the dashed green line.

intriguing properties of overparameterized neural network optimization, it would be appealing (i) to extend our results to deeper network architecture, (ii) to more complex data models, and also (iii) to explore other heuristics that can further boost the robustness of gradient descent methods.

## 6 Proofs

#### 6.1 Proofs for General Theory

We begin by defining the average Jacobian which will be used throughout our analysis.

**Definition 6.1 (Average Jacobian)** We define the average Jacobian along the path connecting two points  $x, y \in \mathbb{R}^p$  as

$$\mathcal{J}(\boldsymbol{y}, \boldsymbol{x}) \coloneqq \int_0^1 \mathcal{J}(\boldsymbol{x} + \alpha(\boldsymbol{y} - \boldsymbol{x})) d\alpha. \tag{6.1}$$

Lemma 6.2 (Linearization of the residual) Given gradient descent iterate  $\hat{\theta} = \theta - \eta \nabla \mathcal{L}(\theta)$ , define

$$C(\theta) = \mathcal{J}(\hat{\theta}, \theta) \mathcal{J}(\theta)^T$$
.

The residuals  $\hat{r} = f(\hat{\theta}) - y$ ,  $r = f(\theta) - y$  obey the following equation

$$\hat{\boldsymbol{r}} = (\boldsymbol{I} - \eta \boldsymbol{C}(\boldsymbol{\theta}))\boldsymbol{r}.$$

**Proof** Following Definition 6.1, denoting  $f(\hat{\theta}) - y = \hat{r}$  and  $f(\theta) - y = r$ , we find that

$$\hat{r} = r - f(\theta) + f(\hat{\theta})$$

$$\stackrel{(a)}{=} r + \mathcal{J}(\hat{\theta}, \theta)(\hat{\theta} - \theta)$$

$$\stackrel{(b)}{=} r - \eta \mathcal{J}(\hat{\theta}, \theta) \mathcal{J}(\theta)^{T} r$$

$$= (I - \eta C(\theta)) r. \tag{6.2}$$

Here (a) uses the fact that Jacobian is the derivative of f and (b) uses the fact that  $\nabla \mathcal{L}(\theta) = \mathcal{J}(\theta)^T r$ .

Using Assumption 3.1, one can show that sparse vectors have small projection on  $\mathcal{S}_+$ .

**Lemma 6.3** Suppose Assumption 3.1 holds. If  $r \in \mathbb{R}^n$  is a vector with s nonzero entries, we have that

$$\|\Pi_{\mathcal{S}_+}(\boldsymbol{r})\|_{\ell_{\infty}} \le \frac{\gamma\sqrt{s}}{n} \|\boldsymbol{r}\|_{\ell_2}.$$
(6.3)

**Proof** First, we bound the  $\ell_2$  projection of r on  $S_+$  as follows

$$\|\Pi_{\mathcal{S}_+}(\boldsymbol{r})\|_{\ell_2} = \sup_{\boldsymbol{v} \in \mathcal{S}_+} \frac{\boldsymbol{v}^T \boldsymbol{r}}{\|\boldsymbol{v}\|_{\ell_2}} \leq \sqrt{\frac{\gamma}{n}} \|\boldsymbol{r}\|_{\ell_1} \leq \sqrt{\frac{\gamma s}{n}} \|\boldsymbol{r}\|_{\ell_2}.$$

where we used the fact that  $|v_i| \le \sqrt{\gamma} ||v||_{\ell_2} / \sqrt{n}$ . Next, we conclude with

$$\|\Pi_{\mathcal{S}_+}(\boldsymbol{r})\|_{\ell_\infty} \leq \sqrt{\frac{\gamma}{n}} \|\Pi_{\mathcal{S}_+}(\boldsymbol{r})\|_{\ell_2} \leq \frac{\gamma\sqrt{s}}{n} \|\boldsymbol{r}\|_{\ell_2}.$$

#### 6.1.1 Proof of Theorem 3.2

**Proof** The proof will be done inductively over the properties of gradient descent iterates and is inspired from the recent work [39]. In particular, [39] requires a well-conditioned Jacobian to fit labels perfectly. In contrast, we have a low-rank Jacobian model which cannot fit the noisy labels (or it would have trouble fitting if the Jacobian was approximately low-rank). Despite this, we wish to prove that gradient descent satisfies desirable properties such as robustness and closeness to initialization. Let us introduce the notation related to the residual. Set  $\mathbf{r}_{\tau} = f(\boldsymbol{\theta}_{\tau}) - \mathbf{y}$  and let  $\mathbf{r}_{0} = f(\boldsymbol{\theta}_{0}) - \mathbf{y}$  be the initial residual. We keep track of the growth of the residual by partitioning the residual as  $\mathbf{r}_{\tau} = \bar{\mathbf{r}}_{\tau} + \bar{\mathbf{e}}_{\tau}$  where

$$ar{e}_{ au} = \Pi_{\mathcal{S}_{-}}(r_{ au})$$
 ,  $ar{r}_{ au} = \Pi_{\mathcal{S}_{+}}(r_{ au})$ .

We claim that for all iterations  $\tau \geq 0$ , the following conditions hold.

$$\bar{\boldsymbol{e}}_{\tau} = \bar{\boldsymbol{e}}_{0} \tag{6.4}$$

$$\|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}^{2} \le \left(1 - \frac{\eta \alpha^{2}}{2}\right)^{\tau} \|\bar{\boldsymbol{r}}_{0}\|_{\ell_{2}}^{2},$$
 (6.5)

$$\frac{1}{4}\alpha \|\boldsymbol{\theta}_{\tau} - \boldsymbol{\theta}_{0}\|_{\ell_{2}} + \|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}} \leq \|\bar{\boldsymbol{r}}_{0}\|_{\ell_{2}} \leq \|\boldsymbol{r}_{0}\|_{\ell_{2}}. \tag{6.6}$$

Assuming these conditions hold till some  $\tau > 0$ , inductively, we focus on iteration  $\tau + 1$ . First, note that these conditions imply that for all  $\tau \ge i \ge 0$ ,  $\theta_i \in \mathcal{D}$  where  $\mathcal{D}$  is the Euclidian ball around  $\theta_0$  of radius  $\frac{4\|\mathbf{r}_0\|_{\ell_2}}{\alpha}$ . This directly follows from (6.6) induction hypothesis. Next, we claim that  $\theta_{\tau+1}$  is still within the set  $\mathcal{D}$ . This can be seen as follows:

Claim 1 Under the induction hypothesis (6.4),  $\theta_{\tau+1} \in \mathcal{D}$ .

**Proof** Since range space of Jacobian is in  $S_+$  and  $\eta \leq 1/\beta^2$ , we begin by noting that

$$\|\boldsymbol{\theta}_{\tau+1} - \boldsymbol{\theta}_{\tau}\|_{\ell_2} = \eta \|\mathcal{J}^T(\boldsymbol{\theta}_{\tau}) \left(f(\boldsymbol{\theta}_{\tau}) - \boldsymbol{y}\right)\|_{\ell_2}$$

$$(6.7)$$

$$\stackrel{(a)}{=} \eta \| \mathcal{J}^T(\boldsymbol{\theta}_\tau) \left( \Pi_{\mathcal{S}_+} (f(\boldsymbol{\theta}_\tau) - \boldsymbol{y}) \right) \|_{\ell_2}$$
(6.8)

$$\stackrel{(b)}{=} \eta \| \mathcal{J}^T(\boldsymbol{\theta}_\tau) \bar{\boldsymbol{r}}_\tau \|_{\ell_2} \tag{6.9}$$

$$\stackrel{(c)}{\leq} \eta \beta \|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2} \tag{6.10}$$

$$\stackrel{(d)}{\leq} \frac{\|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2}}{\beta} \tag{6.11}$$

$$\stackrel{(e)}{\leq} \frac{\|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2}}{\alpha} \tag{6.12}$$

In the above, (a) follows from the fact that row range space of Jacobian is subset of  $S_+$  via Assumption 2. (b) follows from the definition of  $\bar{r}_{\tau}$ . (c) follows from the upper bound on the spectral norm of the Jacobian over  $\mathcal{D}$  per Assumption 2, (d) from the fact that  $\eta \leq \frac{1}{\beta^2}$ , (e) from  $\alpha \leq \beta$ . The latter combined with the triangular inequality and induction hypothesis (6.6) yields (after scaling (6.6) by  $4/\alpha$ )

$$\|\boldsymbol{\theta}_{\tau+1} - \boldsymbol{\theta}_0\|_{\ell_2} \le \|\boldsymbol{\theta}_{\tau+1} - \boldsymbol{\theta}_\tau\|_{\ell_2} + \|\boldsymbol{\theta}_0 - \boldsymbol{\theta}_\tau\|_{\ell_2} \le \|\boldsymbol{\theta}_\tau - \boldsymbol{\theta}_0\|_{\ell_2} + \frac{\|\bar{\boldsymbol{r}}_\tau\|_{\ell_2}}{\alpha} \le \frac{4\|\boldsymbol{r}_0\|_{\ell_2}}{\alpha},$$

concluding the proof of  $\theta_{\tau+1} \in \mathcal{D}$ .

To proceed, we shall verify that (6.6) holds for  $\tau + 1$  as well. Note that, following Lemma 6.2, gradient descent iterate can be written as

$$r_{\tau+1} = (I - C(\theta_{\tau}))r_{\tau}.$$

Since both column and row space of  $C(\theta_{\tau})$  is subset of  $S_{+}$ , we have that

$$\bar{e}_{\tau+1} = \prod_{\mathcal{S}_{-}} ((\boldsymbol{I} - \boldsymbol{C}(\boldsymbol{\theta}_{\tau})) \boldsymbol{r}_{\tau}) \tag{6.13}$$

$$=\Pi_{\mathcal{S}_{-}}(\boldsymbol{r}_{\tau})\tag{6.14}$$

$$= \bar{\boldsymbol{e}}_{\tau}, \tag{6.15}$$

This shows the first statement of the induction. Next, over  $S_+$ , we have

$$\bar{r}_{\tau+1} = \Pi_{\mathcal{S}_+}((I - C(\theta_\tau))r_\tau) \tag{6.16}$$

$$= \Pi_{\mathcal{S}_{+}}((I - C(\theta_{\tau}))\bar{r}_{\tau}) + \Pi_{\mathcal{S}_{+}}((I - C(\theta_{\tau}))\bar{e}_{\tau})$$

$$(6.17)$$

$$= \prod_{S_{\perp}} ((\boldsymbol{I} - \boldsymbol{C}(\boldsymbol{\theta}_{\tau})) \bar{\boldsymbol{r}}_{\tau}) \tag{6.18}$$

$$= (I - C(\theta_{\tau}))\bar{r}_{\tau} \tag{6.19}$$

where the second line uses the fact that  $\bar{e}_{\tau} \in \mathcal{S}_{-}$  and last line uses the fact that  $\bar{r}_{\tau} \in \mathcal{S}_{+}$ . To proceed, we need to prove that  $C(\theta_{\tau})$  has desirable properties over  $\mathcal{S}_{+}$ , in particular, it contracts this space.

Claim 2 let  $P_{S_+} \in \mathbb{R}^{n \times n}$  be the projection matrix to  $S_+$  i.e. it is a positive semi-definite matrix whose eigenvectors over  $S_+$  is 1 and its complement is 0. Under the induction hypothesis and setup of the theorem, we have that<sup>7</sup>

$$\beta^{2} \mathbf{P}_{\mathcal{S}_{+}} \geq \mathbf{C}(\boldsymbol{\theta}_{\tau}) \geq \frac{1}{2} \mathcal{J}(\boldsymbol{\theta}_{\tau}) \mathcal{J}(\boldsymbol{\theta}_{\tau})^{T} \geq \frac{\alpha^{2}}{2} \mathbf{P}_{\mathcal{S}_{+}}.$$
 (6.20)

**Proof** The proof utilizes the upper bound on the learning rate. The argument is similar to the proof of Lemma 9.7 of [39]. Suppose Assumption 3 holds. Then, for any  $\theta_1, \theta_2 \in \mathcal{D}$  we have

$$\|\mathcal{J}(\boldsymbol{\theta}_{2},\boldsymbol{\theta}_{1}) - \mathcal{J}(\boldsymbol{\theta}_{1})\| = \left\| \int_{0}^{1} \left( \mathcal{J}(\boldsymbol{\theta}_{1} + t(\boldsymbol{\theta}_{2} - \boldsymbol{\theta}_{1})) - \mathcal{J}(\boldsymbol{\theta}_{1}) \right) dt \right\|,$$

$$\leq \int_{0}^{1} \|\mathcal{J}(\boldsymbol{\theta}_{1} + t(\boldsymbol{\theta}_{2} - \boldsymbol{\theta}_{1})) - \mathcal{J}(\boldsymbol{\theta}_{1})\| dt,$$

$$\leq \int_{0}^{1} tL \|\boldsymbol{\theta}_{2} - \boldsymbol{\theta}_{1}\|_{\ell_{2}} dt \leq \frac{L}{2} \|\boldsymbol{\theta}_{2} - \boldsymbol{\theta}_{1}\|_{\ell_{2}}.$$

$$(6.21)$$

Thus, for  $\eta \leq \frac{\alpha}{L\beta \|\boldsymbol{r}_0\|_{\ell_2}}$ ,

$$\|\mathcal{J}(\boldsymbol{\theta}_{\tau+1}, \boldsymbol{\theta}_{\tau}) - \mathcal{J}(\boldsymbol{\theta}_{\tau})\| \leq \frac{L}{2} \|\boldsymbol{\theta}_{\tau+1} - \boldsymbol{\theta}_{\tau}\|_{\ell_{2}}$$

$$(6.22)$$

$$= \frac{\eta L}{2} \left\| \mathcal{J}^T(\boldsymbol{\theta}_{\tau}) \left( f(\boldsymbol{\theta}_{\tau}) - \boldsymbol{y} \right) \right\|_{\ell_2} \le \frac{\eta \beta L}{2} \left\| \bar{\boldsymbol{r}}_{\tau} \right\|_{\ell_2}$$
 (6.23)

$$\stackrel{(a)}{\leq} \frac{\eta \beta L}{2} \| \bar{\boldsymbol{r}}_0 \|_{\ell_2} \stackrel{(b)}{\leq} \frac{\alpha}{2}. \tag{6.24}$$

<sup>&</sup>lt;sup>7</sup>We say  $A \ge B$  if A - B is a positive semi-definite matrix in the sense that for any real vector v,  $v^T(A - B)v \ge 0$ .

where for (a) we utilized the induction hypothesis (6.6) and (b) follows from the upper bound on  $\eta$ . Now that (6.24) is established, using following lemma, we find

$$\mathcal{C}(\boldsymbol{\theta}_{\tau}) = \mathcal{J}(\boldsymbol{\theta}_{\tau+1}, \boldsymbol{\theta}_{\tau}) \mathcal{J}(\boldsymbol{\theta}_{\tau})^T \geq (1/2) \mathcal{J}(\boldsymbol{\theta}_{\tau}) \mathcal{J}(\boldsymbol{\theta}_{\tau})^T.$$

The  $\beta^2$  upper bound directly follows from Assumption 2 by again noticing range space of Jacobian is subset of  $S_+$ .

Lemma 6.4 (Asymmetric PSD perturbation) Consider the matrices  $A, C \in \mathbb{R}^{n \times p}$  obeying  $||A - C|| \le \alpha/2$ . Also suppose  $CC^T \ge \alpha^2 P_{S_+}$ . Furthermore, assume range spaces of A, C lies in  $S_+$ . Then,

$$AC^T \ge \frac{CC^T}{2} \ge \frac{\alpha^2}{2} P_{S_+}.$$

**Proof** For  $r \in \mathcal{S}_+$  with unit Euclidian norm, we have

$$\mathbf{r}^{T} \mathbf{A} \mathbf{C}^{T} \mathbf{r} = \|\mathbf{C}^{T} \mathbf{r}\|_{\ell_{2}}^{2} + \mathbf{r}^{T} (\mathbf{A} - \mathbf{C}) \mathbf{C}^{T} \mathbf{r} \ge \|\mathbf{C}^{T} \mathbf{r}\|_{\ell_{2}}^{2} - \|\mathbf{C}^{T} \mathbf{r}\|_{\ell_{2}} \|\mathbf{r}^{T} (\mathbf{A} - \mathbf{C})\|_{\ell_{2}} \\
= (\|\mathbf{C}^{T} \mathbf{r}\|_{\ell_{2}} - \|\mathbf{r}^{T} (\mathbf{A} - \mathbf{C})\|_{\ell_{2}}) \|\mathbf{C}^{T} \mathbf{r}\|_{\ell_{2}} \\
\ge (\|\mathbf{C}^{T} \mathbf{r}\|_{\ell_{2}} - \alpha/2) \|\mathbf{C}^{T} \mathbf{r}\|_{\ell_{2}} \\
\ge \|\mathbf{C}^{T} \mathbf{r}\|_{\ell_{2}}^{2} / 2.$$

Also, for any r, by range space assumption  $r^T A C^T r = \Pi_{\mathcal{S}_+}(r)^T A C^T \Pi_{\mathcal{S}_+}(r)$  (same for  $CC^T$ ). Combined with above, this concludes the claim.

What remains is proving the final two statements of the induction (6.6). Note that, using the claim above and recalling (6.19) and using the fact that  $\|\mathcal{J}(\theta_{\tau+1}, \theta_{\tau})\| \leq \beta$ , the residual satisfies

$$\|\bar{r}_{\tau+1}\|_{\ell_{2}}^{2} = \|(I - \eta C(\theta_{\tau}))\bar{r}_{\tau}\|_{\ell_{2}}^{2} = \|\bar{r}_{\tau}\|_{\ell_{2}}^{2} - 2\eta \bar{r}_{\tau}^{T} C_{\tau} \bar{r}_{\tau} + \eta^{2} \bar{r}_{\tau}^{T} C_{\tau}^{T} C_{\tau} \bar{r}_{\tau}$$

$$(6.25)$$

$$\leq \|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}^{2} - \eta \bar{\boldsymbol{r}}_{\tau}^{T} \mathcal{J}(\boldsymbol{\theta}_{\tau}) \mathcal{J}(\boldsymbol{\theta}_{\tau})^{T} \bar{\boldsymbol{r}}_{\tau} + \eta^{2} \beta^{2} \bar{\boldsymbol{r}}_{\tau}^{T} \mathcal{J}(\boldsymbol{\theta}_{\tau}) \mathcal{J}(\boldsymbol{\theta}_{\tau})^{T} \bar{\boldsymbol{r}}_{\tau}$$

$$(6.26)$$

$$\leq \|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}^{2} - (\eta - \eta^{2}\beta^{2})\|\mathcal{J}(\boldsymbol{\theta}_{\tau})^{T}\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}^{2}$$

$$(6.27)$$

$$\leq \|\bar{r}_{\tau}\|_{\ell_{2}}^{2} - \frac{\eta}{2} \|\mathcal{J}(\boldsymbol{\theta}_{\tau})^{T} \bar{r}_{\tau}\|_{\ell_{2}}^{2}. \tag{6.28}$$

where we used the fact that  $\eta \leq \frac{1}{2\beta^2}$ . Now, using the fact that  $\mathcal{J}(\boldsymbol{\theta}_{\tau})\mathcal{J}(\boldsymbol{\theta}_{\tau})^T \geq \alpha^2 \boldsymbol{P}_{\mathcal{S}_+}$ , we have

$$\|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}^{2} - \frac{\eta}{2} \|\mathcal{J}(\boldsymbol{\theta}_{\tau})^{T} \bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}^{2} \leq (1 - \frac{\eta \alpha^{2}}{2}) \|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}^{2} \leq (1 - \frac{\eta \alpha^{2}}{2})^{\tau+1} \|\bar{\boldsymbol{r}}_{0}\|_{\ell_{2}}^{2},$$

which establishes the second statement of the induction (6.6). What remains is obtaining the last statement of (6.6). To address this, completing squares, observe that

$$\|\bar{\boldsymbol{r}}_{\tau+1}\|_{\ell_2} \leq \sqrt{\|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2}^2 - \frac{\eta}{2} \|\mathcal{J}(\boldsymbol{\theta}_{\tau})^T \bar{\boldsymbol{r}}_{\tau}\|_{\ell_2}^2} \leq \|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2} - \frac{\eta}{4} \frac{\|\mathcal{J}(\boldsymbol{\theta}_{\tau})^T \bar{\boldsymbol{r}}_{\tau}\|_{\ell_2}^2}{\|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2}}.$$

On the other hand, the distance to initial point satisfies

$$\|\boldsymbol{\theta}_{\tau+1} - \boldsymbol{\theta}_0\|_{\ell_2} \le \|\boldsymbol{\theta}_{\tau+1} - \boldsymbol{\theta}_{\tau}\|_{\ell_2} + \|\boldsymbol{\theta}_{\tau} - \boldsymbol{\theta}_0\|_{\ell_2} \le \|\boldsymbol{\theta}_{\tau} - \boldsymbol{\theta}_0\|_{\ell_2} + \eta \|\mathcal{J}(\boldsymbol{\theta}_{\tau})\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2}.$$

Combining the last two lines (by scaling the second line by  $\frac{1}{4}\alpha$ ) and using induction hypothesis (6.6), we find that

$$\frac{1}{4}\alpha \|\boldsymbol{\theta}_{\tau+1} - \boldsymbol{\theta}_0\|_{\ell_2} + \|\bar{\boldsymbol{r}}_{\tau+1}\|_{\ell_2} \le \frac{1}{4}\alpha (\|\boldsymbol{\theta}_{\tau} - \boldsymbol{\theta}_0\|_{\ell_2} + \eta \|\mathcal{J}(\boldsymbol{\theta}_{\tau})\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2}) + \|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2} - \frac{\eta}{4} \frac{\|\mathcal{J}(\boldsymbol{\theta}_{\tau})^T\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2}^2}{\|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2}}$$
(6.29)

$$\leq \left[\frac{1}{4}\alpha \|\boldsymbol{\theta}_{\tau} - \boldsymbol{\theta}_{0}\|_{\ell_{2}} + \|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}\right] + \frac{\eta}{4} \left[\alpha \|\mathcal{J}(\boldsymbol{\theta}_{\tau})\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}} - \frac{\|\mathcal{J}(\boldsymbol{\theta}_{\tau})^{T}\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}^{2}}{\|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}}\right]$$
(6.30)

$$\leq \left[\frac{1}{4}\alpha \|\boldsymbol{\theta}_{\tau} - \boldsymbol{\theta}_{0}\|_{\ell_{2}} + \|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}\right] + \frac{\eta}{4} \|\mathcal{J}(\boldsymbol{\theta}_{\tau})\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}} \left[\alpha - \frac{\|\mathcal{J}(\boldsymbol{\theta}_{\tau})^{T}\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}}{\|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}}}\right]$$
(6.31)

$$\leq \frac{1}{4}\alpha \|\boldsymbol{\theta}_{\tau} - \boldsymbol{\theta}_{0}\|_{\ell_{2}} + \|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{2}} \tag{6.32}$$

$$\leq \|\bar{r}_0\|_{\ell_2} \leq \|r_0\|_{\ell_2}. \tag{6.33}$$

This establishes the final line of the induction and concludes the proof of the upper bound on  $\|\boldsymbol{\theta}_{\tau} - \boldsymbol{\theta}_{0}\|_{\ell_{2}}$ . To proceed, we shall bound the infinity norm of the residual. Using  $\Pi_{\mathcal{S}_{-}}(\boldsymbol{e}) = \Pi_{\mathcal{S}_{-}}(\boldsymbol{r}_{0}) = \bar{\boldsymbol{e}}_{\tau}$ , note that

$$||f(\boldsymbol{\theta}_{\tau}) - \boldsymbol{y} - \boldsymbol{e}||_{\ell_{\infty}} = ||\boldsymbol{r}_{\tau} - \boldsymbol{e}||_{\ell_{\infty}}$$

$$(6.34)$$

$$\leq \|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{\infty}} + \|\boldsymbol{e} - \bar{\boldsymbol{e}}_{\tau}\|_{\ell_{\infty}} \tag{6.35}$$

$$= \|\bar{\mathbf{r}}_{\tau}\|_{\ell_{m}} + \|\mathbf{e} - \Pi_{\mathcal{S}_{-}}(\mathbf{e})\|_{\ell_{m}} \tag{6.36}$$

$$= \|\bar{r}_{\tau}\|_{\ell_{\infty}} + \|\Pi_{\mathcal{S}_{+}}(e)\|_{\ell_{\infty}}. \tag{6.37}$$

What remains is controlling  $\|\bar{r}_{\tau}\|_{\ell_{\infty}}$ . For this term, we shall use the naive upper bound  $\|\bar{r}_{\tau}\|_{\ell_{2}}$ . Using the rate of convergence of the algorithm (6.6), we have that

$$\|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2} \leq (1 - \frac{\eta \alpha^2}{4})^{\tau} \|\boldsymbol{r}_0\|_{\ell_2}.$$

We wish the right hand side to be at most  $\nu > 0$  where  $\nu \ge \|\Pi_{\mathcal{S}_+}(e)\|_{\ell_\infty}$ . This implies that we need

$$(1 - \frac{\eta \alpha^2}{4})^{\tau} \| \mathbf{r}_0 \|_{\ell_2} \le \nu \iff \tau \log(1 - \frac{\eta \alpha^2}{4}) \le \log(\frac{\nu}{\| \mathbf{r}_0 \|_{\ell_2}})$$
(6.38)

$$\iff \tau \log\left(\frac{1}{1 - \frac{\eta \alpha^2}{4}}\right) \ge \log\left(\frac{\|r_0\|_{\ell_2}}{\nu}\right) \tag{6.39}$$

To conclude, note that since  $\frac{\eta\alpha^2}{4} \le 1/8$  (as  $\eta \le 1/2\beta^2$ ), we have

$$\log \left(\frac{1}{1 - \frac{\eta \alpha^2}{4}}\right) \ge \log \left(1 + \frac{\eta \alpha^2}{4}\right) \ge \frac{\eta \alpha^2}{5}.$$

Consequently, if  $\tau \geq \frac{5}{\eta \alpha^2} \log(\frac{\|\boldsymbol{r}_0\|_{\ell_2}}{\nu})$ , we find that  $\|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_{\infty}} \leq \|\bar{\boldsymbol{r}}_{\tau}\|_{\ell_2} \leq \nu$ , which guarantees

$$\|\boldsymbol{r}_{\tau} - \boldsymbol{e}\|_{\ell_{\infty}} \leq 2\nu.$$

which is the advertised result. If e is s sparse and  $S_+$  is diffused, applying Lemma 3.1 we have

$$\|\Pi_{\mathcal{S}_+}(\boldsymbol{e})\|_{\ell_\infty} \leq \frac{\gamma\sqrt{s}}{n} \|\boldsymbol{e}\|_{\ell_2}.$$

#### 6.1.2 Proof of Generic Lower Bound - Theorem 3.3

**Proof** Suppose  $\theta \in \mathcal{D}$  satisfies  $\mathbf{y} = f(\theta)$ . Define  $\mathbf{J}_{\tau} = \mathcal{J}((1-\tau)\theta + \tau\theta_0)$  and  $\mathbf{J} = \mathcal{J}(\theta, \theta_0) = \int_0^1 \mathbf{J}_{\tau} d\tau$ . Since Jacobian is derivative of f, we have that

$$f(\boldsymbol{\theta}) - f(\boldsymbol{\theta}_0) = \int_0^1 \boldsymbol{J}_{\tau}(\boldsymbol{\theta} - \boldsymbol{\theta}_0) d\tau = \boldsymbol{J}(\boldsymbol{\theta} - \boldsymbol{\theta}_0).$$

Now, define the matrices  $J_+ = \Pi_{S_+}(J)$  and  $J_- = \Pi_{S_-}(J)$ . Using Assumption 1, we bound the spectral norms via

$$\|\boldsymbol{J}_{+}\| = \sup_{\boldsymbol{v} \in \mathcal{S}_{+}, \|\boldsymbol{v}\|_{\ell_{2}} \leq 1} \|\boldsymbol{J}^{T}\boldsymbol{v}\|_{\ell_{2}} \leq \beta \quad , \quad \|\boldsymbol{J}_{-}\| = \sup_{\boldsymbol{v} \in \mathcal{S}_{-}, \|\boldsymbol{v}\|_{\ell_{2}} \leq 1} \|\boldsymbol{J}^{T}\boldsymbol{v}\|_{\ell_{2}} \leq \epsilon.$$

To proceed, projecting the residual on  $S_+$ , we find for any  $\theta$  with  $f(\theta) = y$ 

$$\Pi_{\mathcal{S}_+}(f(\boldsymbol{\theta}) - f(\boldsymbol{\theta}_0)) = \Pi_{\mathcal{S}_+}(\boldsymbol{J})(\boldsymbol{\theta} - \boldsymbol{\theta}_0) \implies \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_{\ell_2} \ge \frac{\|\Pi_{\mathcal{S}_+}(f(\boldsymbol{\theta}) - f(\boldsymbol{\theta}_0))\|_{\ell_2}}{\beta} \ge \frac{E_+}{\beta}.$$

The identical argument for  $S_{-}$  yields  $\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_{\ell_2} \geq \frac{E_{-}}{\epsilon}$ . Together this implies

$$\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_{\ell_2} \ge \max(\frac{E_-}{\epsilon}, \frac{E_+}{\beta}). \tag{6.40}$$

If R is strictly smaller than right hand side, we reach a contradiction as  $\theta \notin \mathcal{D}$ . If  $\mathcal{D} = \mathbb{R}^p$ , we still find (6.40).

This shows that if  $\epsilon$  is small and  $E_{-}$  is nonzero, gradient descent has to traverse a long distance to find a good model. Intuitively, if the projection over the noise space indeed contains the label noise, we actually don't want to fit that. Algorithmically, our idea fits the residual over the signal space and not worries about fitting over the noise space. Approximately speaking, this intuition corresponds to the  $\ell_2$  regularized problem

$$\min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) \qquad \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_{\ell_2} \leq R.$$

If we set  $R = \frac{E_+}{\beta}$ , we can hope that solution will learn only the signal and does not overfit to the noise. The next section builds on this intuition and formalizes our algorithmic guarantees.

#### 6.2 Proofs for Neural Networks

Throughout,  $\sigma_{\min}(\cdot)$  denotes the smallest singular value of a given matrix. We first introduce helpful definitions that will be used in our proofs.

**Definition 6.5 (Support subspace)** Let  $\{x_i\}_{i=1}^n$  be an input dataset generated according to Definition 1.1. Also let  $\{\tilde{x}_i\}_{i=1}^n$  be the associated cluster centers, that is,  $\tilde{x}_i = \mathbf{c}_\ell$  iff  $\mathbf{x}_i$  is from the  $\ell$ th cluster. We define the support subspace  $S_+$  as a subspace of dimension K, dictated by the cluster membership as follows. Let  $\Lambda_\ell \subset \{1, \ldots, n\}$  be the set of coordinates i such that  $\tilde{x}_i = \mathbf{c}_\ell$ . Then,  $S_+$  is characterized by

$$\mathcal{S}_+ = \{ \boldsymbol{v} \in \mathbb{R}^n \mid \boldsymbol{v}_{i_1} = \boldsymbol{v}_{i_2} \quad \textit{for all} \quad i_1, i_2 \in \Lambda_\ell \quad \textit{and for all } 1 \leq \ell \leq K \}.$$

**Definition 6.6 (Neural Net Jacobian)** Given input samples  $(\boldsymbol{x}_i)_{i=1}^n$ , form the input matrix  $\boldsymbol{X} = [\boldsymbol{x}_1 \dots \boldsymbol{x}_n]^T \in \mathbb{R}^{n \times d}$ . The Jacobian of the learning problem (1.3), at a matrix  $\boldsymbol{W}$  is denoted by  $\mathcal{J}(\boldsymbol{W}, \boldsymbol{X}) \in \mathbb{R}^{n \times kd}$  and is given by

$$\mathcal{J}(\boldsymbol{W},\boldsymbol{X})^T = (diag(\boldsymbol{v})\phi'(\boldsymbol{W}\boldsymbol{X}^T)) * \boldsymbol{X}^T.$$

Here \* denotes the Khatri-Rao product.

The following theorem is borrowed from [40] and characterizes three key properties of the neural network Jacobian. These are smoothness, spectral norm, and minimum singular value at initialization which correspond to Lemmas 6.6, 6.7, and 6.8 in that paper.

Theorem 6.7 (Jacobian Properties at Cluster Center) Suppose  $X = [x_1 \dots x_n]^T \in \mathbb{R}^{n \times d}$  be an input dataset satisfying  $\lambda(X) > 0$ . Suppose  $|\phi'|, |\phi''| \leq \Gamma$ . The Jacobian mapping with respect to the input-to-hidden weights obey the following properties.

• Smoothness is bounded by

$$\|\mathcal{J}(\widetilde{\boldsymbol{W}}, \boldsymbol{X}) - \mathcal{J}(\boldsymbol{W}, \boldsymbol{X})\| \le \frac{\Gamma}{\sqrt{k}} \|\boldsymbol{X}\| \|\widetilde{\boldsymbol{W}} - \boldsymbol{W}\|_F \quad for \ all \quad \widetilde{\boldsymbol{W}}, \boldsymbol{W} \in \mathbb{R}^{k \times d}.$$

• Top singular value is bounded by

$$\|\mathcal{J}(\boldsymbol{W},\boldsymbol{X})\| \leq \Gamma \|\boldsymbol{X}\|$$
.

• Let C > 0 be an absolute constant. As long as

$$k \ge \frac{C\Gamma^2 \log n \|\boldsymbol{X}\|^2}{\lambda(\boldsymbol{X})}$$

At random Gaussian initialization  $W_0 \sim \mathcal{N}(0,1)^{k \times d}$ , with probability at least  $1 - 1/K^{100}$ , we have

$$\sigma_{\min}\left(\mathcal{J}(\boldsymbol{W}_0, \boldsymbol{X})\right) \geq \sqrt{\lambda(\boldsymbol{X})/2}.$$

In our case, the Jacobian is **not** well-conditioned. However, it is pretty well-structured as described previously. To proceed, given a matrix  $X \in \mathbb{R}^{n \times d}$  and a subspace  $S \subset \mathbb{R}^n$ , we define the minimum singular value of the matrix over this subspace by  $\sigma_{\min}(X, S)$  which is defined as

$$\sigma_{\min}(\boldsymbol{X}, \mathcal{S}) = \sup_{\|\boldsymbol{v}\|_{\ell_2} = 1, \boldsymbol{U}\boldsymbol{U}^T = \boldsymbol{P}_{\mathcal{S}}} \|\boldsymbol{v}^T \boldsymbol{U}^T \boldsymbol{X}\|_{\ell_2}.$$

Here,  $P_{\mathcal{S}} \in \mathbb{R}^{n \times n}$  is the projection operator to the subspace. Hence, this definition essentially projects the matrix on  $\mathcal{S}$  and then takes the minimum singular value over that projected subspace. The following theorem states the properties of the Jacobian at a clusterable dataset.

**Theorem 6.8 (Jacobian Properties at Clusterable Dataset)** Let input samples  $(x_i)_{i=1}^n$  be generated according to  $(\varepsilon_0, \delta)$  clusterable dataset model of Definition 1.1 and define  $X = [x_1 \dots x_n]^T$ . Let  $S_+$  be the support space and  $(\tilde{x}_i)_{i=1}^n$  be the associated clean dataset as described by Definition 6.5. Set  $\tilde{X} = [\tilde{x}_1 \dots \tilde{x}_n]^T$ . Assume  $|\phi'|, |\phi''| \leq \Gamma$  and  $\lambda(C) > 0$ . The Jacobian mapping at  $\tilde{X}$  with respect to the input-to-hidden weights obey the following properties.

• Smoothness is bounded by

$$\left\| \mathcal{J}(\widetilde{\boldsymbol{W}}, \widetilde{\boldsymbol{X}}) - \mathcal{J}(\boldsymbol{W}, \widetilde{\boldsymbol{X}}) \right\| \leq \Gamma \sqrt{\frac{c_{up}n}{kK}} \left\| \boldsymbol{C} \right\| \left\| \widetilde{\boldsymbol{W}} - \boldsymbol{W} \right\|_{F} \quad \textit{for all} \quad \widetilde{\boldsymbol{W}}, \boldsymbol{W} \in \mathbb{R}^{k \times d}.$$

• Top singular value is bounded by

$$\|\mathcal{J}(\boldsymbol{W}, \tilde{\boldsymbol{X}})\| \leq \sqrt{\frac{c_{up}n}{K}} \Gamma \|\boldsymbol{C}\|.$$

• As long as

$$k \ge \frac{C\Gamma^2 \log K \|\boldsymbol{C}\|^2}{\lambda(\boldsymbol{C})}$$

At random Gaussian initialization  $W_0 \sim \mathcal{N}(0,1)^{k \times d}$ , with probability at least  $1 - 1/K^{100}$ , we have

$$\sigma_{\min}\left(\mathcal{J}(\boldsymbol{W}_0, \tilde{\boldsymbol{X}}), \mathcal{S}_+\right) \ge \sqrt{\frac{c_{low}n\lambda(\boldsymbol{C})}{2K}}$$

• The range space obeys range  $(\mathcal{J}(W_0, \tilde{X})) \subset \mathcal{S}_+$  where  $\mathcal{S}_+$  is given by Definition 6.5.

**Proof** Let  $\mathcal{J}(W,C)$  be the Jacobian at the cluster center matrix. Applying Theorem 6.7, this matrix already obeys the properties described in the conclusions of this theorem with desired probability (for the last conclusion). We prove our theorem by relating the cluster center Jacobian to the clean dataset Jacobian matrix  $\mathcal{J}(W,\tilde{X})$ .

Note that  $\tilde{X}$  is obtained by duplicating the rows of the cluster center matrix C. This implies that  $\mathcal{J}(W, \tilde{X})$  is obtained by duplicating the rows of the cluster center Jacobian. The critical observation is that, by construction in Definition 1.1, each row is duplicated somewhere between  $c_{low} n/K$  and  $c_{up} n/K$ .

To proceed, fix a vector v and let  $\tilde{p} = \mathcal{J}(W, \tilde{X})v \in \mathbb{R}^n$  and  $p = \mathcal{J}(W, C)v \in \mathbb{R}^K$ . Recall the definition of the support sets  $\Lambda_{\ell}$  from Definition 6.5. We have the identity

$$\tilde{\boldsymbol{p}}_i = \boldsymbol{p}_\ell$$
 for all  $i \in \Lambda_\ell$ .

This implies  $\tilde{p} \in \mathcal{S}_+$  hence range $(\mathcal{J}(W, \tilde{X})) \subset \mathcal{S}_+$ . Furthermore, the entries of  $\tilde{p}$  repeats the entries of p somewhere between  $c_{low}n/K$  and  $c_{up}n/K$ . This implies that,

$$\sqrt{\frac{c_{up}n}{K}} \|\boldsymbol{p}\|_{\ell_2} \geq \|\tilde{\boldsymbol{p}}\|_{\ell_2} \geq \sqrt{\frac{c_{low}n}{K}} \|\boldsymbol{p}\|_{\ell_2},$$

and establishes the upper and lower bounds on the singular values of  $\mathcal{J}(W, \tilde{X})$  over  $\mathcal{S}_+$  in terms of the singular values of  $\mathcal{J}(W, C)$ . Finally, the smoothness can be established similarly. Given matrices  $W, \tilde{W}$ , the rows of the difference

$$\|\mathcal{J}(\widetilde{oldsymbol{W}}, ilde{oldsymbol{X}}) - \mathcal{J}(oldsymbol{W}, ilde{oldsymbol{X}})\|$$

is obtained by duplicating the rows of  $\|\mathcal{J}(\widetilde{W}, C) - \mathcal{J}(W, C)\|$  by at most  $c_{up}n/K$  times. Hence the spectral norm is scaled by at most  $\sqrt{c_{up}n/K}$ .

Lemma 6.9 (Upper bound on initial misfit) Consider a one-hidden layer neural network model of the form  $\mathbf{x} \mapsto \mathbf{v}^T \phi(\mathbf{W}\mathbf{x})$  where the activation  $\phi$  has bounded derivatives obeying  $|\phi(0)|, |\phi'(z)| \leq \Gamma$ . Suppose entries of  $\mathbf{v} \in \mathbb{R}^k$  are half  $1/\sqrt{k}$  and half  $-1/\sqrt{k}$  so that  $\|\mathbf{v}\|_{\ell_2} = 1$ . Also assume we have n data points  $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n \in \mathbb{R}^d$  with unit euclidean norm ( $\|\mathbf{x}_i\|_{\ell_2} = 1$ ) aggregated as rows of a matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and the corresponding labels given by  $\mathbf{y} \in \mathbb{R}^n$  generated according to  $(\rho, \varepsilon_0 = 0, \delta)$  noisy dataset (Definition 1.2). Then for  $\mathbf{W}_0 \in \mathbb{R}^{k \times d}$  with i.i.d.  $\mathcal{N}(0, 1)$  entries

$$\|\boldsymbol{v}^T\phi\left(\boldsymbol{W}_0\boldsymbol{X}^T\right)-\boldsymbol{y}\|_{\ell_2} \leq \mathcal{O}(\Gamma\sqrt{n\log K}),$$

holds with probability at least  $1 - K^{-100}$ .

**Proof** This lemma is based on a fairly straightforward union bound. First, by construction  $\|\boldsymbol{y}\|_{\ell_2} \leq \sqrt{n}$ . What remains is bounding  $\|\boldsymbol{v}^T\phi(\boldsymbol{W}_0\boldsymbol{X}^T)\|_{\ell_2}$ . Since  $\varepsilon_0 = 0$  there are K unique rows. We will show that each of the unique rows is bounded with probability  $1 - K^{-101}$  and union bounding will give the final result. Let  $\boldsymbol{w}$  be a row of  $\boldsymbol{W}_0$  and  $\boldsymbol{x}$  be a row of  $\boldsymbol{X}$ . Since  $\phi$  is  $\Gamma$  Lipschitz and  $|\phi(0)| \leq \Gamma$ , each entry of  $\phi(\boldsymbol{X}\boldsymbol{w})$  is  $\mathcal{O}(\Gamma)$ -subgaussian. Hence  $\boldsymbol{v}^T\phi(\boldsymbol{W}_0\boldsymbol{x})$  is weighted average of k i.i.d. subgaussians which are entries of  $\phi(\boldsymbol{W}_0\boldsymbol{x})$ . Additionally it is zero mean since  $\sum_{i=1}^n \boldsymbol{v}_i = 0$ . This means  $\boldsymbol{v}^T\phi(\boldsymbol{W}_0\boldsymbol{x})$  is also  $\mathcal{O}(\Gamma)$  subgaussian and obeys

$$\mathbb{P}(|\boldsymbol{v}^T \phi(\boldsymbol{W}_0 \boldsymbol{x})| \ge c\Gamma \sqrt{\log K}) \le K^{-101},$$

for some constant c > 0, concluding the proof.

#### 6.2.1 Proof of Theorem 2.3

We first prove a lemma regarding the projection of label noise on the cluster induced subspace.

**Lemma 6.10** Let  $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$  be an  $(\rho, \varepsilon_0 = 0, \delta)$  clusterable noisy dataset as described in Definition 1.2. Let  $\{\tilde{y}_i\}_{i=1}^n$  be the corresponding noiseless labels. Let  $\mathcal{J}(\boldsymbol{W}, \boldsymbol{C})$  be the Jacobian at the cluster center matrix which is rank K and  $\mathcal{S}_+$  be its column space. Then, the difference between noiseless and noisy labels satisfy the bound

$$\|\Pi_{\mathcal{S}_{\perp}}(\boldsymbol{y}-\tilde{\boldsymbol{y}})\|_{\ell_{\infty}} \leq 2\rho.$$

**Proof** Let  $e = y - \tilde{y}$ . Observe that by assumption,  $\ell$ th cluster has at most  $s_{\ell} = \rho n_{\ell}$  errors. Let  $\mathcal{I}_{\ell}$  denote the membership associated with cluster  $\ell$  i.e.  $\mathcal{I}_{\ell} \subset \{1, \ldots, n\}$  and  $i \in \mathcal{I}_{\ell}$  if and only if  $x_i$  belongs to  $\ell$ th cluster. Let  $\mathbf{1}(\ell) \in \mathbb{R}^n$  be the indicator function of the  $\ell$ th class where ith entry is 1 if  $i \in \mathcal{I}_{\ell}$  and 0 else for  $1 \le i \le n$ . Then, denoting the size of the  $\ell$ th cluster by  $n_{\ell}$ , the projection to subspace  $\mathcal{S}_{+}$  can be written as the P matrix where

$$\boldsymbol{P} = \sum_{\ell=1}^{K} \frac{1}{n_{\ell}} \mathbf{1}(\ell) \mathbf{1}(\ell)^{T}.$$

Let  $e_{\ell}$  be the error pattern associated with  $\ell$ th cluster i.e.  $e_{\ell}$  is equal to e over  $\mathcal{I}_{\ell}$  and zero outside. Since cluster membership is non-overlapping, we have that

$$Pe = \sum_{\ell=1}^{K} \frac{1}{n_{\ell}} \mathbf{1}(\ell) \mathbf{1}(\ell)^{T} e_{\ell}.$$

Similarly since supports of  $\mathbf{1}(\ell)$  are non-overlapping, we have that

$$\|\boldsymbol{P}\boldsymbol{e}\|_{\ell_{\infty}} = \max_{1 \le \ell \le K} \frac{1}{n_{\ell}} \mathbf{1}(\ell) \mathbf{1}(\ell)^{T} \boldsymbol{e}_{\ell}.$$

Now, using  $\|e\|_{\ell_{\infty}} \le 2$  (max distance between two labels), observe that

$$\|\mathbf{1}(\ell)\mathbf{1}(\ell)^T e_{\ell}\|_{\ell_{\infty}} \le 2\|\mathbf{1}(\ell)\|_{\ell_{\infty}} \|e_{\ell}\|_{\ell_{1}} = 2\|e_{\ell}\|_{\ell_{1}}.$$

Since number of errors within cluster  $\ell$  is at most  $n_{\ell}\rho$ , we find that

$$\|Pe\|_{\ell_{\infty}} = \sum_{\ell=1}^{K} \|\frac{1}{n_{\ell}} \mathbf{1}(\ell) \mathbf{1}(\ell)^{T} e_{\ell}\|_{\ell_{\infty}} \le \frac{\|e_{\ell}\|_{\ell_{1}}}{n_{\ell}} \le 2\rho.$$

The final line yields the bound

$$\|\mathcal{P}_{S_{+}}(y-\tilde{y})\|_{\ell_{\infty}} = \|\mathcal{P}_{S_{+}}(e)\|_{\ell_{\infty}} = \|Pe\|_{\ell_{\infty}} \le 2\rho.$$

With this, we are ready to state the proof of Theorem 2.3.

**Proof** The proof is based on the meta Theorem 3.2, hence we need to verify its Assumptions 2 and 3 with proper values and apply Lemma 6.10 to get  $\|\mathcal{P}_{\mathcal{S}_+}(e)\|_{\ell_{\infty}}$ . We will also make significant use of Corollary 6.8.

Using Corollary 6.8, Assumption 3 holds with  $L = \Gamma \sqrt{\frac{c_{up}n}{kK}} \| \boldsymbol{C} \|$  where L is the Lipschitz constant of Jacobian spectrum. Denote  $\boldsymbol{r}_{\tau} = f(\boldsymbol{W}_{\tau}) - \boldsymbol{y}$ . Using Lemma 6.9 with probability  $1 - K^{-100}$ , we have that  $\|\boldsymbol{r}_0\|_{\ell_2} = \|\boldsymbol{y} - f(\boldsymbol{W}_0)\|_{\ell_2} \le \Gamma \sqrt{c_0 n \log K/128}$  for some  $c_0 > 0$ . Corollary 6.8 guarantees a uniform bound for  $\beta$ , hence in Assumption 2, we pick

$$\beta \leq \sqrt{\frac{c_{up}n}{K}} \Gamma \| \boldsymbol{C} \|.$$

We shall also pick the minimum singular value over  $\mathcal{S}_+$  to be

$$\alpha = \frac{\alpha_0}{2}$$
 where  $\alpha_0 = \sqrt{\frac{c_{low}n\lambda(C)}{2K}}$ ,

We wish to verify Assumption 2 over the radius of

$$R = \frac{4\|f(\boldsymbol{W}_0) - \boldsymbol{y}\|_{\ell_2}}{\alpha} \le \frac{\Gamma\sqrt{c_0 n \log K/8}}{\alpha} = \Gamma\sqrt{\frac{c_0 n \log K/2}{\frac{c_{low} n \lambda(\boldsymbol{C})}{2K}}} = \Gamma\sqrt{\frac{c_0 K \log K}{c_{low} \lambda(\boldsymbol{C})}},$$

neighborhood of  $W_0$ . What remains is ensuring that Jacobian over  $S_+$  is lower bounded by  $\alpha$ . Our choice of k guarantees that at the initialization, with probability  $1 - K^{-100}$ , we have

$$\sigma_{\min}(\mathcal{J}(\boldsymbol{W}_0, \boldsymbol{X}), \mathcal{S}_+) \geq \alpha_0.$$

Suppose  $LR \leq \alpha = \alpha_0/2$ . Using triangle inequality on Jacobian spectrum, for any  $\mathbf{W} \in \mathcal{D}$ , using  $\|\mathbf{W} - \mathbf{W}_0\|_F \leq R$ , we would have

$$\sigma_{\min}(\mathcal{J}(\boldsymbol{W}, \boldsymbol{X}), \mathcal{S}_+) \geq \sigma_{\min}(\mathcal{J}(\boldsymbol{W}_0, \boldsymbol{X}), \mathcal{S}_+) - LR \geq \alpha_0 - \alpha = \alpha.$$

Now, observe that

$$LR = \Gamma \sqrt{\frac{c_{up}n}{kK}} \| \boldsymbol{C} \| \Gamma \sqrt{\frac{c_0 K \log(K)}{c_{low} \lambda(\boldsymbol{C})}} = \Gamma^2 \| \boldsymbol{C} \| \sqrt{\frac{c_{up}c_0 n \log K}{c_{low} k \lambda(\boldsymbol{C})}} \le \frac{\alpha_0}{2} = \sqrt{\frac{c_{low}n\lambda(\boldsymbol{C})}{8K}}, \quad (6.41)$$

as k satisfies

$$k \ge \mathcal{O}(\Gamma^4 \| \boldsymbol{C} \|^2 \frac{c_{up} K \log(K)}{c_{low}^2 \lambda(\boldsymbol{C})^2}) \ge \mathcal{O}(\frac{\Gamma^4 K \log(K) \| \boldsymbol{C} \|^2}{\lambda(\boldsymbol{C})^2}).$$

Finally, since  $LR = 4L \| \mathbf{r}_0 \|_{\ell_2} / \alpha \le \alpha$ , the learning rate is

$$\eta \leq \frac{1}{2\beta^2} \min(1, \frac{\alpha\beta}{L \|\boldsymbol{r}_0\|_{\ell_2}}) = \frac{1}{2\beta^2} = \frac{K}{2c_{up}n\Gamma^2 \|\boldsymbol{C}\|^2}.$$

Overall, the assumptions of Theorem 3.2 holds with stated  $\alpha, \beta, L$  with probability  $1-2K^{-100}$  (union bounding initial residual and minimum singular value events). This implies for all  $\tau > 0$  the distance of current iterate to initial obeys

$$\|\boldsymbol{W}_{\tau} - \boldsymbol{W}_0\|_F \leq R.$$

The final step is the properties of the label corruption. Using Lemma 6.10, we find that

$$\|\Pi_{\mathcal{S}_{+}}(\tilde{\boldsymbol{y}}-\boldsymbol{y})\|_{\ell_{++}} \leq 2\rho.$$

Substituting the values corresponding to  $\alpha, \beta, L$  yields that, for all gradient iterations with

$$\frac{5}{\eta \alpha^2} \log(\frac{\|\boldsymbol{r}_0\|_{\ell_2}}{2\rho}) \le \frac{5}{\eta \alpha^2} \log(\frac{\Gamma \sqrt{c_0 n \log K/32}}{2\rho}) = \mathcal{O}(\frac{K}{\eta n \lambda(\boldsymbol{C})} \log(\frac{\Gamma \sqrt{n \log K}}{\rho})) \le \tau,$$

denoting the clean labels by  $\tilde{y}$  and applying Theorem 3.2, we have that, the infinity norm of the residual obeys (using  $\|\Pi_{\mathcal{S}_+}(e)\|_{\ell_\infty} \leq 2\rho$ )

$$||f(\boldsymbol{W}) - \tilde{\boldsymbol{y}}||_{\ell_{\infty}} \le 4\rho.$$

This implies that if  $\rho \le \delta/8$ , the network will miss the correct label by at most  $\delta/2$ , hence all labels (including noisy ones) will be correctly classified.

#### 6.2.2 Proof of Theorem 2.4

Consider

$$f(\boldsymbol{W}, \boldsymbol{x}) = \boldsymbol{v}^T \phi(\boldsymbol{W} \boldsymbol{x})$$

and note that

$$\nabla_{\boldsymbol{x}} f(\boldsymbol{W}, \boldsymbol{x}) = \boldsymbol{W}^T \operatorname{diag}(\phi'(\boldsymbol{W}\boldsymbol{x})) \boldsymbol{v}$$

Thus

$$\frac{\partial}{\partial x} f(\boldsymbol{W}, \boldsymbol{x}) \boldsymbol{u} = \boldsymbol{v}^T \operatorname{diag} \left( \phi'(\boldsymbol{W} \boldsymbol{x}) \right) \boldsymbol{W} \boldsymbol{u}$$
$$= \sum_{\ell=1}^k \boldsymbol{v}_\ell \phi'(\langle \boldsymbol{w}_\ell, \boldsymbol{x} \rangle) \boldsymbol{w}_\ell^T \boldsymbol{u}$$

Thus

$$\nabla_{\boldsymbol{w}_{\ell}} \left( \frac{\partial}{\partial \boldsymbol{x}} f(\boldsymbol{W}, \boldsymbol{x}) \boldsymbol{u} \right) = \boldsymbol{v}_{\ell} \left( \phi''(\boldsymbol{w}_{\ell}^T \boldsymbol{x}) (\boldsymbol{w}_{\ell}^T \boldsymbol{u}) \boldsymbol{x} + \phi'(\boldsymbol{w}_{\ell}^T \boldsymbol{x}) \boldsymbol{u} \right)$$

Thus, denoting vectorization of a matrix by  $\text{vect}(\cdot)$ 

$$\begin{aligned} \operatorname{vect}(\boldsymbol{U})^T \left( \frac{\partial}{\partial \operatorname{vect}(\boldsymbol{W})} \frac{\partial}{\partial \boldsymbol{x}} f(\boldsymbol{W}, \boldsymbol{x}) \right) \boldsymbol{u} &= \sum_{\ell=1}^k \boldsymbol{v}_\ell \left( \phi''(\boldsymbol{w}_\ell^T \boldsymbol{x}) (\boldsymbol{w}_\ell^T \boldsymbol{u}) (\boldsymbol{u}_\ell^T \boldsymbol{x}) + \phi'(\boldsymbol{w}_\ell^T \boldsymbol{x}) (\boldsymbol{u}_\ell^T \boldsymbol{u}) \right) \\ &= \boldsymbol{u}^T \boldsymbol{W}^T \operatorname{diag}(\boldsymbol{v}) \operatorname{diag}(\boldsymbol{\phi}''(\boldsymbol{W} \boldsymbol{x})) \boldsymbol{U} \boldsymbol{x} + \boldsymbol{v}^T \operatorname{diag}(\boldsymbol{\phi}'(\boldsymbol{W} \boldsymbol{x})) \boldsymbol{U} \boldsymbol{u} \end{aligned}$$

Thus by the general mean value theorem there exists a point  $(\widetilde{W}, \widetilde{x})$  in the square  $(W_0, x_1), (W_0, x_2), (W, x_1)$  and  $(W, x_2)$  such that

$$(f(\boldsymbol{W}, \boldsymbol{x}_2) - f(\boldsymbol{W}_0, \boldsymbol{x}_2)) - (f(\boldsymbol{W}, \boldsymbol{x}_1) - f(\boldsymbol{W}_0, \boldsymbol{x}_1))$$

$$= (\boldsymbol{x}_2 - \boldsymbol{x}_1)^T \widetilde{\boldsymbol{W}}^T \operatorname{diag}(\boldsymbol{v}) \operatorname{diag}(\phi''(\widetilde{\boldsymbol{W}}\widetilde{\boldsymbol{x}})) (\boldsymbol{W} - \boldsymbol{W}_0) \widetilde{\boldsymbol{x}} + \boldsymbol{v}^T \operatorname{diag}(\phi'(\widetilde{\boldsymbol{W}}\widetilde{\boldsymbol{x}})) (\boldsymbol{W} - \boldsymbol{W}_0) (\boldsymbol{x}_2 - \boldsymbol{x}_1)$$

Using the above we have that

$$\left| (f(\boldsymbol{W}, \boldsymbol{x}_{2}) - f(\boldsymbol{W}_{0}, \boldsymbol{x}_{2})) - (f(\boldsymbol{W}, \boldsymbol{x}_{1}) - f(\boldsymbol{W}_{0}, \boldsymbol{x}_{1})) \right|$$

$$\stackrel{(a)}{\leq} \left| (\boldsymbol{x}_{2} - \boldsymbol{x}_{1})^{T} \widetilde{\boldsymbol{W}}^{T} \operatorname{diag}(\boldsymbol{v}) \operatorname{diag}\left(\phi''(\widetilde{\boldsymbol{W}}\widetilde{\boldsymbol{x}})\right) (\boldsymbol{W} - \boldsymbol{W}_{0}) \widetilde{\boldsymbol{x}} \right|$$

$$+ \left| \boldsymbol{v}^{T} \operatorname{diag}\left(\phi'\left(\widetilde{\boldsymbol{W}}\widetilde{\boldsymbol{x}}\right)\right) (\boldsymbol{W} - \boldsymbol{W}_{0}) (\boldsymbol{x}_{2} - \boldsymbol{x}_{1}) \right|$$

$$\stackrel{(b)}{\leq} \left( \left\| \boldsymbol{v} \right\|_{\ell_{\infty}} \left\| \widetilde{\boldsymbol{x}} \right\|_{\ell_{2}} \left\| \widetilde{\boldsymbol{W}} \right\| + \left\| \boldsymbol{v} \right\|_{\ell_{2}} \right) \Gamma \left\| \boldsymbol{x}_{2} - \boldsymbol{x}_{1} \right\|_{\ell_{2}} \left\| \boldsymbol{W} - \boldsymbol{W}_{0} \right\|$$

$$\stackrel{(c)}{\leq} \left( \frac{1}{\sqrt{k}} \left\| \widetilde{\boldsymbol{W}} \right\| + 1 \right) \Gamma \left\| \boldsymbol{x}_{2} - \boldsymbol{x}_{1} \right\|_{\ell_{2}} \left\| \boldsymbol{W} - \boldsymbol{W}_{0} \right\|$$

$$\stackrel{(e)}{\leq} \left( \frac{1}{\sqrt{k}} \left\| \boldsymbol{W}_{0} \right\| + \frac{1}{\sqrt{k}} \left\| \widetilde{\boldsymbol{W}} - \boldsymbol{W}_{0} \right\| + 1 \right) \Gamma \left\| \boldsymbol{x}_{2} - \boldsymbol{x}_{1} \right\|_{\ell_{2}} \left\| \boldsymbol{W} - \boldsymbol{W}_{0} \right\|$$

$$\stackrel{(f)}{\leq} \left( \frac{1}{\sqrt{k}} \left\| \boldsymbol{W}_{0} \right\| + \frac{1}{\sqrt{k}} \left\| \widetilde{\boldsymbol{W}} - \boldsymbol{W}_{0} \right\|_{F} + 1 \right) \Gamma \left\| \boldsymbol{x}_{2} - \boldsymbol{x}_{1} \right\|_{\ell_{2}} \left\| \boldsymbol{W} - \boldsymbol{W}_{0} \right\|$$

$$\stackrel{(g)}{\leq} \left( \frac{1}{\sqrt{k}} \left\| \widetilde{\boldsymbol{W}} - \boldsymbol{W}_{0} \right\|_{F} + 3 + 2\sqrt{\frac{d}{k}} \right) \Gamma \left\| \boldsymbol{x}_{2} - \boldsymbol{x}_{1} \right\|_{\ell_{2}} \left\| \boldsymbol{W} - \boldsymbol{W}_{0} \right\|$$

$$\stackrel{(h)}{\leq} C \Gamma \left\| \boldsymbol{x}_{2} - \boldsymbol{x}_{1} \right\|_{\ell_{2}} \left\| \boldsymbol{W} - \boldsymbol{W}_{0} \right\|$$

$$\stackrel{(h)}{\leq} C \Gamma \left\| \boldsymbol{x}_{2} - \boldsymbol{x}_{1} \right\|_{\ell_{2}} \left\| \boldsymbol{W} - \boldsymbol{W}_{0} \right\|$$

$$\stackrel{(h)}{\leq} C \Gamma \left\| \boldsymbol{x}_{2} - \boldsymbol{x}_{1} \right\|_{\ell_{2}} \left\| \boldsymbol{W} - \boldsymbol{W}_{0} \right\|$$

$$(6.42)$$

Here, (a) follows from the triangle inequality, (b) from simple algebraic manipulations along with the fact that  $|\phi'(z)| \le \Gamma$  and  $|\phi''(z)| \le \Gamma$ , (c) from the fact that  $\mathbf{v}_{\ell} = \pm \frac{1}{\sqrt{k}}$ , (d) from  $\|\mathbf{x}_2\|_{\ell_2} = \|\mathbf{x}_1\|_{\ell_2} = 1$  which implies  $\|\widetilde{\mathbf{x}}\|_{\ell_2} \le 1$ , (e) from triangular inequality, (f) from the fact that Frobenius norm dominates the spectral norm, (g) from the fact that with probability at least  $1 - 2e^{-(d+k)}$ ,  $\|\mathbf{W}_0\| \le 2(\sqrt{k} + \sqrt{d})$ , and (h) from the fact that  $\|\widetilde{\mathbf{W}} - \mathbf{W}_0\| \le \|\mathbf{W} - \mathbf{W}_0\|_F \le \widetilde{c}\sqrt{k}$  and  $k \ge cd$ .

Next we note that for a Gaussian random vector  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  we have

$$\|\phi(\boldsymbol{g}^{T}\boldsymbol{x}_{2}) - \phi(\boldsymbol{g}^{T}\boldsymbol{x}_{1})\|_{\psi_{2}} = \|\phi(\boldsymbol{g}^{T}\boldsymbol{x}_{2}) - \phi(\boldsymbol{g}^{T}\boldsymbol{x}_{1})\|_{\psi_{2}}$$

$$= \|\phi'\left(t\boldsymbol{g}^{T}\boldsymbol{x}_{2} + (1-t)\boldsymbol{g}^{T}\boldsymbol{x}_{1}\right)\boldsymbol{g}^{T}(\boldsymbol{x}_{2} - \boldsymbol{x}_{1})\|_{\psi_{2}}$$

$$\leq \Gamma \|\boldsymbol{g}^{T}(\boldsymbol{x}_{2} - \boldsymbol{x}_{1})\|_{\psi_{2}}$$

$$\leq c\Gamma \|\boldsymbol{x}_{2} - \boldsymbol{x}_{1}\|_{\ell_{2}}.$$
(6.43)

Also note that

$$egin{aligned} f(oldsymbol{W}_0, oldsymbol{x}_2) - f(oldsymbol{W}_0, oldsymbol{x}_1) = & oldsymbol{v}^T \left( \phi\left(oldsymbol{W}_0 oldsymbol{x}_2 
ight) - \phi\left(oldsymbol{W}_0 oldsymbol{x}_1 
ight) 
ight) \ & \sim \sum_{\ell=1}^k oldsymbol{v}_\ell \left( \phi(oldsymbol{g}_\ell^T oldsymbol{x}_2) - \phi(oldsymbol{g}_\ell^T oldsymbol{x}_1) 
ight) \end{aligned}$$

where  $g_1, g_2, ..., g_k$  are i.i.d. vectors with  $\mathcal{N}(0, I_d)$  distribution. Also for  $\mathbf{v}$  obeying  $\mathbf{1}^T \mathbf{v} = 0$  this random variable has mean zero. Hence, using the fact that weighted sum of subGaussian random variables are subgaussian combined with (B.2) we conclude that  $f(\mathbf{W}_0, \mathbf{x}_2) - f(\mathbf{W}_0, \mathbf{x}_1)$  is also subGaussian obeying  $||f(\mathbf{W}_0, \mathbf{x}_2) - f(\mathbf{W}_0, \mathbf{x}_1)||_{\psi_2} \le c\Gamma ||\mathbf{v}||_{\ell_2} ||\mathbf{x}_2 - \mathbf{x}_1||_{\ell_2}$ . Thus

$$|f(\mathbf{W}_{0}, \mathbf{x}_{2}) - f(\mathbf{W}_{0}, \mathbf{x}_{1})| \le ct\Gamma \|\mathbf{v}\|_{\ell_{2}} \|\mathbf{x}_{2} - \mathbf{x}_{1}\|_{\ell_{2}} = ct\Gamma \|\mathbf{x}_{2} - \mathbf{x}_{1}\|_{\ell_{2}},$$
 (6.44)

with probability at least  $1 - e^{-\frac{t^2}{2}}$ .

Now combining (B.1) and (B.3) we have

$$\delta \leq |y_{2} - y_{2}|$$

$$= |f(\boldsymbol{W}, \boldsymbol{x}_{1}) - f(\boldsymbol{W}, \boldsymbol{x}_{2})|$$

$$= |\boldsymbol{v}^{T} (\phi(\boldsymbol{W}\boldsymbol{x}_{2}) - \phi(\boldsymbol{W}\boldsymbol{x}_{1}))|$$

$$\leq |(f(\boldsymbol{W}, \boldsymbol{x}_{2}) - f(\boldsymbol{W}_{0}, \boldsymbol{x}_{2})) - (f(\boldsymbol{W}, \boldsymbol{x}_{1}) - f(\boldsymbol{W}_{0}, \boldsymbol{x}_{1}))| + |\boldsymbol{v}^{T} (\phi(\boldsymbol{W}_{0}\boldsymbol{x}_{2}) - \phi(\boldsymbol{W}_{0}\boldsymbol{x}_{1}))|$$

$$\leq C\Gamma \|\boldsymbol{x}_{2} - \boldsymbol{x}_{1}\|_{\ell_{2}} \|\boldsymbol{W} - \boldsymbol{W}_{0}\| + ct\Gamma \|\boldsymbol{x}_{2} - \boldsymbol{x}_{1}\|_{\ell_{2}}$$

$$\leq C\Gamma \varepsilon_{0} \left( \|\boldsymbol{W} - \boldsymbol{W}_{0}\| + \frac{1}{1000}t \right)$$

Thus

$$\|\boldsymbol{W} - \boldsymbol{W}_0\| \ge \frac{\delta}{C\Gamma\varepsilon_0} - \frac{t}{1000},$$

with high probability.

#### 6.3 Perturbation analysis for perfectly clustered data (Proof of Theorem 2.2)

Denote average neural net Jacobian at data  $\boldsymbol{X}$  via

$$\mathcal{J}(\boldsymbol{W}_1, \boldsymbol{W}_2, \boldsymbol{X}) = \int_0^1 \mathcal{J}(\alpha \boldsymbol{W}_1 + (1 - \alpha) \boldsymbol{W}_2, \boldsymbol{X}) d\alpha.$$

**Lemma 6.11 (Perturbed Jacobian Distance)** Let  $X = [x_1 \dots x_n]^T$  be the input matrix obtained from Definition 1.1. Let  $\tilde{X}$  be the noiseless inputs where  $\tilde{x}_i$  is the cluster center corresponding to  $x_i$ . Given weight matrices  $W_1, W_2, \tilde{W}_1, \tilde{W}_2$ , we have that

$$\|\mathcal{J}(\boldsymbol{W}_{1}, \boldsymbol{W}_{2}, \boldsymbol{X}) - \mathcal{J}(\tilde{\boldsymbol{W}}_{1}, \tilde{\boldsymbol{W}}_{2}, \tilde{\boldsymbol{X}})\| \leq \Gamma \sqrt{n} (\frac{\|\tilde{\boldsymbol{W}}_{1} - \boldsymbol{W}_{1}\|_{F} + \|\tilde{\boldsymbol{W}}_{2} - \boldsymbol{W}_{2}\|_{F}}{2\sqrt{k}} + \varepsilon_{0}).$$

**Proof** Given  $W, \tilde{W}$ , we write

$$\|\mathcal{J}(\boldsymbol{W},\boldsymbol{X}) - \mathcal{J}(\tilde{\boldsymbol{W}},\tilde{\boldsymbol{X}})\| \leq \|\mathcal{J}(\boldsymbol{W},\boldsymbol{X}) - \mathcal{J}(\tilde{\boldsymbol{W}},\boldsymbol{X})\| + \|\mathcal{J}(\tilde{\boldsymbol{W}},\boldsymbol{X}) - \mathcal{J}(\tilde{\boldsymbol{W}},\tilde{\boldsymbol{X}})\|.$$

We first bound

$$\|\mathcal{J}(\boldsymbol{W}, \boldsymbol{X}) - \mathcal{J}(\tilde{\boldsymbol{W}}, \boldsymbol{X})\| = \|\operatorname{diag}(\boldsymbol{v})\phi'(\boldsymbol{W}\boldsymbol{X}^T) * \boldsymbol{X}^T - \operatorname{diag}(\boldsymbol{v})\phi'(\tilde{\boldsymbol{W}}\boldsymbol{X}^T) * \boldsymbol{X}^T\|$$
(6.45)

$$= \frac{1}{\sqrt{k}} \| (\phi'(\boldsymbol{W}\boldsymbol{X}^T) - \phi'(\tilde{\boldsymbol{W}}\boldsymbol{X}^T)) * \boldsymbol{X}^T \|$$
 (6.46)

To proceed, we use the results on the spectrum of Hadamard product of matrices due to Schur [46]. Given  $\mathbf{A} \in \mathbb{R}^{k \times d}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times d}$  matrices where  $\mathbf{B}$  has unit length rows, we have

$$\|\boldsymbol{A} * \boldsymbol{B}\| = \sqrt{\|(\boldsymbol{A} * \boldsymbol{B})^T (\boldsymbol{A} * \boldsymbol{B})\|} = \sqrt{\|(\boldsymbol{A}^T \boldsymbol{A}) \odot (\boldsymbol{B}^T \boldsymbol{B})\|} \le \sqrt{\|\boldsymbol{A}^T \boldsymbol{A}\|} = \|\boldsymbol{A}\|.$$

Substituting  $\mathbf{A} = \phi'(\mathbf{W}\mathbf{X}^T) - \phi'(\tilde{\mathbf{W}}\mathbf{X}^T)$  and  $\mathbf{B} = \mathbf{X}^T$ , we find

$$\|(\phi'(\boldsymbol{W}\boldsymbol{X}^T) - \phi'(\tilde{\boldsymbol{W}}\boldsymbol{X}^T)) * \boldsymbol{X}^T\| \leq \|\phi'(\boldsymbol{W}\boldsymbol{X}^T) - \phi'(\tilde{\boldsymbol{W}}\boldsymbol{X}^T)\| \leq \Gamma \|(\tilde{\boldsymbol{W}} - \boldsymbol{W})\boldsymbol{X}^T\|_F \leq \Gamma \sqrt{n} \|\tilde{\boldsymbol{W}} - \boldsymbol{W}\|_F.$$

Secondly,

$$\|\mathcal{J}(\tilde{\boldsymbol{W}},\boldsymbol{X}) - \mathcal{J}(\tilde{\boldsymbol{W}},\tilde{\boldsymbol{X}})\| = \frac{1}{\sqrt{k}} \|\phi'(\tilde{\boldsymbol{W}}\boldsymbol{X}^T) * (\boldsymbol{X} - \tilde{\boldsymbol{X}})\|$$

where reusing Schur's result and boundedness of  $|\phi'| \leq \Gamma$ 

$$\|\phi'(\tilde{\boldsymbol{W}}\boldsymbol{X}^T)*(\boldsymbol{X}-\tilde{\boldsymbol{X}})\| \leq \Gamma\sqrt{k}\|\boldsymbol{X}-\tilde{\boldsymbol{X}}\| \leq \Gamma\sqrt{kn}\varepsilon_0.$$

Combining both estimates yields

$$\|\mathcal{J}(\boldsymbol{W}, \boldsymbol{X}) - \mathcal{J}(\tilde{\boldsymbol{W}}, \tilde{\boldsymbol{X}})\| \leq \Gamma \sqrt{n} (\frac{\|\tilde{\boldsymbol{W}} - \boldsymbol{W}\|_F}{\sqrt{k}} + \varepsilon_0).$$

To get the result on  $\|\mathcal{J}(W_1, W_2, X) - \mathcal{J}(\tilde{W}_1, \tilde{W}_2, \tilde{X})\|$ , we integrate

$$\|\mathcal{J}(\boldsymbol{W}_{1}, \boldsymbol{W}_{2}, \boldsymbol{X}) - \mathcal{J}(\tilde{\boldsymbol{W}}_{1}, \tilde{\boldsymbol{W}}_{2}, \tilde{\boldsymbol{X}})\| \leq \int_{0}^{1} \Gamma \sqrt{n} \left(\frac{\|\alpha(\tilde{\boldsymbol{W}}_{1} - \boldsymbol{W}_{1}) + (1 - \alpha)(\tilde{\boldsymbol{W}}_{1} - \boldsymbol{W}_{1})\|_{F}}{\sqrt{k}} + \varepsilon_{0}\right) d\alpha \qquad (6.47)$$

$$\leq \Gamma \sqrt{n} \left( \frac{\|\tilde{W}_1 - W_1\|_F + \|\tilde{W}_2 - W_2\|_F}{2\sqrt{k}} + \varepsilon_0 \right). \tag{6.48}$$

Theorem 6.12 (Robustness of gradient path to perturbation) Generate samples  $(\boldsymbol{x}_i, y_i)_{i=1}^n$  according to  $(\rho, \varepsilon_0, \delta)$  noisy dataset model and form the concatenated input/labels  $\boldsymbol{X} \in \mathbb{R}^{d \times n}, \boldsymbol{y} \in \mathbb{R}^n$ . Let  $\tilde{\boldsymbol{X}}$  be the clean input sample matrix obtained by mapping  $\boldsymbol{x}_i$  to its associated cluster center. Set learning rate  $\eta \leq \frac{K}{2c_{up}n\Gamma^2\|\boldsymbol{C}\|^2}$  and maximum iterations  $\tau_0$  satisfying

$$\eta \tau_0 = C_1 \frac{K}{n\lambda(C)} \log(\frac{\Gamma \sqrt{n \log K}}{\rho}).$$

where  $C_1 \ge 1$  is a constant of our choice. Suppose input noise level  $\varepsilon_0$  and number of hidden nodes obey

$$\varepsilon_0 \le \mathcal{O}\left(\frac{\lambda(C)}{\Gamma^2 K \log(\frac{\Gamma \sqrt{n \log K}}{\rho})}\right) \quad and \quad k \ge \mathcal{O}\left(\Gamma^{10} \frac{K^2 \|C\|^4}{\lambda(C)^4} \log(\frac{\Gamma \sqrt{n \log K}}{\rho})^6\right).$$

Set  $W_0 \stackrel{i.i.d.}{\sim} \mathcal{N}(0,1)$ . Starting from  $W_0 = \tilde{W}_0$  consider the gradient descent iterations over the losses

$$\mathbf{W}_{\tau+1} = \mathbf{W}_{\tau} - \eta \nabla \mathcal{L}(\mathbf{W}_{\tau}) \quad \text{where} \quad \mathcal{L}(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^{n} (y_i - f(\mathbf{W}, \tilde{\mathbf{x}}_i))^2$$
(6.49)

$$\tilde{\boldsymbol{W}}_{\tau+1} = \tilde{\boldsymbol{W}}_{\tau} - \nabla \tilde{\mathcal{L}}(\tilde{\boldsymbol{W}}_{\tau}) \quad \text{where} \quad \tilde{\mathcal{L}}(\tilde{\boldsymbol{W}}) = \frac{1}{2} \sum_{i=1}^{n} (y_i - f(\tilde{\boldsymbol{W}}, \tilde{\boldsymbol{x}}_i))^2$$
(6.50)

Then, for all gradient descent iterations satisfying  $\tau \leq \tau_0$ , we have that

$$||f(\boldsymbol{W}_{\tau}, \boldsymbol{X}) - f(\tilde{\boldsymbol{W}}_{\tau}, \tilde{\boldsymbol{X}})||_{\ell_2} \le c_0 \tau \eta \varepsilon_0 \Gamma^3 n^{3/2} \sqrt{\log K},$$

and

$$\|\boldsymbol{W}_{\tau} - \tilde{\boldsymbol{W}}_{\tau}\|_{F} \leq \mathcal{O}(\tau \eta \varepsilon_{0} \frac{\Gamma^{4} K n}{\lambda(\boldsymbol{C})} \log(\frac{\Gamma \sqrt{n \log K}}{\rho})^{2}).$$

**Proof** Since  $\tilde{W}_{\tau}$  are the noiseless iterations, with probability  $1 - 2K^{-100}$ , the statements of Theorem 2.3 hold on  $\tilde{W}_{\tau}$ . To proceed with proof, we first introduce short hand notations. We use

$$r_i = f(\mathbf{W}_i, \mathbf{X}) - \mathbf{y}, \ \tilde{r}_i = f(\tilde{\mathbf{W}}_i, \tilde{\mathbf{X}}_i) - \mathbf{y}$$

$$(6.51)$$

$$\mathcal{J}_{i} = \mathcal{J}(\boldsymbol{W}_{i}, \boldsymbol{X}), \ \mathcal{J}_{i+1,i} = \mathcal{J}(\boldsymbol{W}_{i+1}, \boldsymbol{W}_{i}, \boldsymbol{X}), \ \tilde{\mathcal{J}}_{i} = \mathcal{J}(\tilde{\boldsymbol{W}}_{i}, \tilde{\boldsymbol{X}}), \ \tilde{\mathcal{J}}_{i+1,i} = \mathcal{J}(\tilde{\boldsymbol{W}}_{i+1}, \tilde{\boldsymbol{W}}_{i}, \tilde{\boldsymbol{X}})$$
(6.52)

$$d_{i} = \|\boldsymbol{W}_{i} - \tilde{\boldsymbol{W}}_{i}\|_{F}, \ p_{i} = \|\boldsymbol{r}_{i} - \tilde{\boldsymbol{r}}_{i}\|_{F}, \ \beta = \Gamma \|\boldsymbol{C}\| \sqrt{c_{up}n/K}, \ L = \Gamma \|\boldsymbol{C}\| \sqrt{c_{up}n/Kk}.$$

$$(6.53)$$

Here  $\beta$  is the upper bound on the Jacobian spectrum and L is the spectral norm Lipschitz constant as in Theorem 6.8. Applying Lemma 6.11, note that

$$\|\mathcal{J}(\boldsymbol{W}_{\tau}, \boldsymbol{X}) - \mathcal{J}(\tilde{\boldsymbol{W}}_{\tau}, \tilde{\boldsymbol{X}})\| \le L\|\tilde{\boldsymbol{W}} - \boldsymbol{W}\|_{F} + \Gamma\sqrt{n\varepsilon_{0}} \le Ld_{\tau} + \Gamma\sqrt{n\varepsilon_{0}}$$

$$(6.54)$$

$$\|\mathcal{J}(\boldsymbol{W}_{\tau+1}, \boldsymbol{W}_{\tau}, \boldsymbol{X}) - \mathcal{J}(\tilde{\boldsymbol{W}}_{\tau+1}, \tilde{\boldsymbol{W}}_{\tau}, \tilde{\boldsymbol{X}})\| \le L(d_{\tau} + d_{\tau+1})/2 + \Gamma\sqrt{n\varepsilon_0}. \tag{6.55}$$

Following this and using that noiseless residual is non-increasing and satisfies  $\|\tilde{r}_{\tau}\|_{\ell_2} \leq \|\tilde{r}_0\|_{\ell_2}$ , note that parameter satisfies

$$\mathbf{W}_{i+1} = \mathbf{W}_i - \eta \mathcal{J}_i \mathbf{r}_i \quad , \quad \tilde{\mathbf{W}}_{i+1} = \tilde{\mathbf{W}}_i - \eta \tilde{\mathcal{J}}_i^T \tilde{\mathbf{r}}_i$$
 (6.56)

$$\|\boldsymbol{W}_{i+1} - \tilde{\boldsymbol{W}}_{i+1}\|_{F} \leq \|\boldsymbol{W}_{i} - \tilde{\boldsymbol{W}}_{i}\|_{F} + \eta \|\mathcal{J}_{i} - \tilde{\mathcal{J}}_{i}\| \|\tilde{\boldsymbol{r}}_{i}\|_{\ell_{2}} + \eta \|\mathcal{J}_{i}\| \|\boldsymbol{r}_{i} - \tilde{\boldsymbol{r}}_{i}\|_{\ell_{2}}$$

$$(6.57)$$

$$d_{i+1} \le d_i + \eta \left( \left( Ld_i + \Gamma \sqrt{n\varepsilon_0} \right) \|\tilde{r}_0\|_{\ell_2} + \beta p_i \right), \tag{6.58}$$

and residual satisfies (using  $I \geq \tilde{\mathcal{J}}_{i+1,i} \tilde{\mathcal{J}}_i^T / \beta^2 \geq 0$ )

$$\mathbf{r}_{i+1} = \mathbf{r}_i - \eta \mathcal{J}_{i+1} \,_i \mathcal{J}_i^T \mathbf{r}_i \Longrightarrow \tag{6.59}$$

$$\boldsymbol{r}_{i+1} - \tilde{\boldsymbol{r}}_{i+1} = (\boldsymbol{r}_i - \tilde{\boldsymbol{r}}_i) - \eta(\mathcal{J}_{i+1,i} - \tilde{\mathcal{J}}_{i+1,i})\mathcal{J}_i^T \boldsymbol{r}_i - \eta \tilde{\mathcal{J}}_{i+1,i}(\mathcal{J}_i^T - \tilde{\mathcal{J}}_i^T) \boldsymbol{r}_i - \eta \tilde{\mathcal{J}}_{i+1,i} \tilde{\mathcal{J}}_i^T (\boldsymbol{r}_i - \tilde{\boldsymbol{r}}_i). \tag{6.60}$$

$$\boldsymbol{r}_{i+1} - \tilde{\boldsymbol{r}}_{i+1} = (\boldsymbol{I} - \eta \tilde{\mathcal{J}}_{i+1,i} \tilde{\mathcal{J}}_{i}^{T})(\boldsymbol{r}_{i} - \tilde{\boldsymbol{r}}_{i}) - \eta (\mathcal{J}_{i+1,i} - \tilde{\mathcal{J}}_{i+1,i}) \mathcal{J}_{i}^{T} \boldsymbol{r}_{i} - \eta \tilde{\mathcal{J}}_{i+1,i} (\mathcal{J}_{i}^{T} - \tilde{\mathcal{J}}_{i}^{T}) \boldsymbol{r}_{i}.$$
(6.61)

$$\|\mathbf{r}_{i+1} - \tilde{\mathbf{r}}_{i+1}\|_{\ell_2} \le \|\mathbf{r}_i - \tilde{\mathbf{r}}_i\|_{\ell_2} + \eta \beta \|\mathbf{r}_i\|_{\ell_2} (L(3d_{\tau} + d_{\tau+1})/2 + 2\Gamma \sqrt{n\varepsilon_0}).$$
(6.62)

$$\|\mathbf{r}_{i+1} - \tilde{\mathbf{r}}_{i+1}\|_{\ell_2} \le \|\mathbf{r}_i - \tilde{\mathbf{r}}_i\|_{\ell_2} + \eta\beta(\|\tilde{\mathbf{r}}_0\|_{\ell_2} + p_i)(L(3d_{\tau} + d_{\tau+1})/2 + 2\Gamma\sqrt{n\varepsilon_0}). \tag{6.63}$$

where we used  $\|\boldsymbol{r}_i\|_{\ell_2} \leq p_i + \|\tilde{\boldsymbol{r}}_0\|_{\ell_2}$  and  $\|(\boldsymbol{I} - \eta \tilde{\mathcal{J}}_{i+1,i} \tilde{\mathcal{J}}_i^T) \boldsymbol{v}\|_{\ell_2} \leq \|\boldsymbol{v}\|_{\ell_2}$  which follows from (6.28). This implies

$$p_{i+1} \le p_i + \eta \beta(\|\tilde{r}_0\|_{\ell_2} + p_i) (L(3d_{\tau} + d_{\tau+1})/2 + 2\Gamma \sqrt{n\varepsilon_0}). \tag{6.64}$$

Finalizing proof: Next, using Lemma 6.9, we have  $\|\tilde{r}_0\|_{\ell_2} \leq \Theta := C_0 \Gamma \sqrt{n \log K}$ . We claim that if

$$\varepsilon_0 \le \mathcal{O}\left(\frac{1}{\tau_0 \eta \Gamma^2 n}\right) \le \frac{1}{8\tau_0 \eta \beta \Gamma \sqrt{n}} \quad \text{and} \quad L \le \frac{2}{5\tau_0 \eta \Theta(1 + 8\eta \tau_0 \beta^2)} \le \frac{1}{30(\tau_0 \eta \beta)^2 \Theta},$$
(6.65)

(where we used  $\eta \tau_0 \beta^2 \ge 1$ ), for all  $t \le \tau_0$ , we have that

$$p_t \le 8t\eta\Gamma\sqrt{n\varepsilon_0}\Theta\beta \le \Theta$$
 ,  $d_t \le 2t\eta\Gamma\sqrt{n\varepsilon_0}\Theta(1+8\eta\tau_0\beta^2)$ . (6.66)

The proof is by induction. Suppose it holds until  $t \le \tau_0 - 1$ . At t + 1, via (6.58) we have that

$$\frac{d_{t+1} - d_t}{\eta} \le L d_t \Theta + \Gamma \sqrt{n} \varepsilon_0 \Theta + 8\tau_0 \eta \beta^2 \Gamma \sqrt{n} \varepsilon_0 \Theta \stackrel{?}{\le} 2\Gamma \sqrt{n} \varepsilon_0 \Theta (1 + 8\eta \tau_0 \beta^2).$$

Right hand side holds since  $L \leq \frac{1}{2\eta\tau_0\Theta}$ . This establishes the induction for  $d_{t+1}$ .

Next, we show the induction on  $p_t$ . Observe that  $3d_t + d_{t+1} \le 10\tau_0\eta\Gamma\sqrt{n\varepsilon_0\Theta}(1+8\eta\tau_0\beta^2)$ . Following (6.64) and using  $p_t \leq \Theta$ , we need

$$\frac{p_{t+1} - p_t}{\eta} \le \beta \Theta \left( L(3d_\tau + d_{\tau+1}) + 4\Gamma \sqrt{n\varepsilon_0} \right) \stackrel{?}{\le} 8\Gamma \sqrt{n\varepsilon_0} \Theta \beta \iff (6.67)$$

$$L(3d_{\tau} + d_{\tau+1}) + 4\Gamma\sqrt{n\varepsilon_0} \stackrel{?}{\leq} 8\Gamma\sqrt{n\varepsilon_0} \iff (6.68)$$

$$L(3d_{\tau} + d_{\tau+1}) \stackrel{?}{\leq} 4\Gamma \sqrt{n}\varepsilon_0 \iff (6.69)$$

$$10L\tau_0\eta(1+8\eta\tau_0\beta^2)\Theta \stackrel{?}{\leq} 4 \iff (6.70)$$

$$L \stackrel{?}{\leq} \frac{2}{5\tau_0 \eta (1 + 8\eta \tau_0 \beta^2) \Theta}. \tag{6.71}$$

Concluding the induction since L satisfies the final line. Consequently, for all  $0 \le t \le \tau_0$ , we have that

$$p_t \leq 8t\eta \Gamma \sqrt{n\varepsilon_0} \Theta \beta = c_0 t\eta \varepsilon_0 \Gamma^3 n^{3/2} \sqrt{\log K}.$$

Next, note that, condition on L is implied by

$$k \ge 1000\Gamma^2 n(\tau_0 \eta \beta)^4 \Theta^2 \tag{6.72}$$

$$= \mathcal{O}(\Gamma^4 n \frac{K^4}{n^4 \lambda(C)^4} \log(\frac{\Gamma \sqrt{n \log K}}{\rho})^4 (\|C\| \Gamma \sqrt{n/K})^4 (\Gamma \sqrt{n \log K})^2)$$
(6.73)

$$= \mathcal{O}\left(\Gamma^{10} \frac{K^2 \|\boldsymbol{C}\|^4}{\lambda(\boldsymbol{C})^4} \log\left(\frac{\Gamma\sqrt{n\log K}}{\rho}\right)^4 \log^2(K)\right)$$
(6.74)

which is implied by  $k \ge \mathcal{O}(\Gamma^{10} \frac{K^2 \|\mathcal{C}\|^4}{\lambda(\mathcal{C})^4} \log(\frac{\Gamma \sqrt{n \log K}}{\rho})^6)$ . Finally, following (6.66), distance satisfies

$$d_t \le 20t\eta^2 \tau_0 \Gamma \sqrt{n\varepsilon_0} \Theta \beta^2 \le \mathcal{O}(t\eta \varepsilon_0 \frac{\Gamma^4 K n}{\lambda(C)} \log(\frac{\Gamma \sqrt{n\log K}}{\rho})^2).$$

#### Completing the Proof of Theorem 2.2

Theorem 2.2 is obtained by the theorem below when we ignore the log terms, and treating  $\Gamma$ ,  $\lambda(C)$  as constants. We also plug in  $\eta = \frac{K}{2c_{mn}n\Gamma^2 \|C\|^2}$ .

Theorem 6.13 (Training neural nets with corrupted labels) Let  $\{(x_i, y_i)\}_{i=1}^n$  be an  $(s, \varepsilon_0, \delta)$  clusterable noisy dataset as described in Definition 1.2. Let  $\{\tilde{y}_i\}_{i=1}^n$  be the corresponding noiseless labels. Suppose  $|\phi(0)|, |\phi'|, |\phi''| \leq \Gamma$  for some  $\Gamma \geq 1$ , input noise and the number of hidden nodes satisfy

$$\varepsilon_0 \leq \mathcal{O}\left(\frac{\lambda(\boldsymbol{C})}{\Gamma^2 K \log\left(\frac{\Gamma \sqrt{n \log K}}{\rho}\right)}\right) \quad and \quad k \geq \mathcal{O}\left(\Gamma^{10} \frac{K^2 \|\boldsymbol{C}\|^4}{\lambda(\boldsymbol{C})^4} \log\left(\frac{\Gamma \sqrt{n \log K}}{\rho}\right)^6\right).$$

where  $C \in \mathbb{R}^{K \times d}$  is the matrix of cluster centers. Set learning rate  $\eta \leq \frac{K}{2c_{up}n\Gamma^2\|C\|^2}$  and randomly initialize  $W_0 \stackrel{i.i.d.}{\sim} \mathcal{N}(0,1)$ . With probability  $1 - 3/K^{100}$ , after  $\tau = \mathcal{O}(\frac{K}{\eta n \lambda(C)}) \log(\frac{\Gamma \sqrt{n \log K}}{\rho})$  iterations, for all  $1 \leq i \leq n$ , we have that

• The per sample normalized  $\ell_2$  norm bound satisfies

$$\frac{\|f(\boldsymbol{W}_{\tau}, \boldsymbol{X}) - \tilde{\boldsymbol{y}}\|_{\ell_2}}{\sqrt{n}} \le 4\rho + c \frac{\varepsilon_0 \Gamma^3 K \sqrt{\log K}}{\lambda(\boldsymbol{C})} \log(\frac{\Gamma \sqrt{n \log K}}{\rho}).$$

• Suppose  $\rho \leq \delta/8$ . Denote the total number of prediction errors with respect to true labels (i.e. not satisfying (2.2)) by  $err(\mathbf{W})$ . With same probability,  $err(\mathbf{W}_{\tau})$  obeys

$$\frac{err(\boldsymbol{W}_{\tau})}{n} \leq c \frac{\varepsilon_0 K}{\delta} \frac{\Gamma^3 \sqrt{\log K}}{\lambda(\boldsymbol{C})} \log(\frac{\Gamma \sqrt{n \log K}}{\rho}).$$

- Suppose  $\rho \leq \delta/8$  and  $\varepsilon_0 \leq c' \frac{\delta \lambda(C)^2}{\Gamma^5 K^2 \log(\frac{\Gamma \sqrt{n \log K}}{\rho})^3}$ , then,  $\mathbf{W}_{\tau}$  assigns all input samples  $\mathbf{x}_i$  to correct ground truth labels  $\tilde{y}_i$  i.e. (2.2) holds for all  $1 \leq i \leq n$ .
- Finally, for any iteration count  $0 \le t \le \tau$  the total distance to initialization is bounded as

$$\|\boldsymbol{W}_{\tau} - \boldsymbol{W}_{0}\|_{F} \leq \mathcal{O}\left(\Gamma\sqrt{\frac{K\log K}{\lambda(\boldsymbol{C})}} + t\eta\varepsilon_{0}\frac{\Gamma^{4}Kn}{\lambda(\boldsymbol{C})}\log\left(\frac{\Gamma\sqrt{n\log K}}{\rho}\right)^{2}\right). \tag{6.75}$$

**Proof** Note that proposed number of iterations  $\tau$  is set so that it is large enough for Theorem 2.3 to achieve small error in the clean input model ( $\varepsilon_0 = 0$ ) and it is small enough so that Theorem 6.12 is applicable. In light of Theorems 6.12 and 2.3 consider two gradient descent iterations starting from  $W_0$  where one uses clean dataset (as if input vectors are perfectly cluster centers)  $\tilde{X}$  and other uses the original dataset X. Denote the prediction residual vectors of the noiseless and original problems at time  $\tau$  with respect true ground truth labels  $\tilde{y}$  by  $\tilde{r}_{\tau} = f(\tilde{W}_{\tau}, \tilde{X}) - \tilde{y}$  and  $r_{\tau} = f(W_{\tau}, X) - \tilde{y}$  respectively. Applying Theorems 6.12 and 2.3, under the stated conditions, we have that

$$\|\tilde{\boldsymbol{r}}_{\tau}\|_{\ell_{\infty}} \le 4\rho \quad \text{and}$$
 (6.76)

$$\|\boldsymbol{r}_{\tau} - \tilde{\boldsymbol{r}}_{\tau}\|_{\ell_{2}} \le c\varepsilon_{0} \frac{K}{n\lambda(\boldsymbol{C})} \log(\frac{\Gamma\sqrt{n\log K}}{\rho}) \Gamma^{3} n^{3/2} \sqrt{\log K}$$

$$(6.77)$$

$$= c \frac{\varepsilon_0 \Gamma^3 K \sqrt{n \log K}}{\lambda(C)} \log(\frac{\Gamma \sqrt{n \log K}}{\rho})$$
 (6.78)

First statement: The latter two results imply the  $\ell_2$  error bounds on  $r_{\tau} = f(\boldsymbol{W}_{\tau}, \boldsymbol{X}) - \tilde{\boldsymbol{y}}$ . Second statement: To assess the classification rate we count the number of entries of  $r_{\tau} = f(\boldsymbol{W}_{\tau}, \boldsymbol{X}) - \tilde{\boldsymbol{y}}$  that is larger than the class margin  $\delta/2$  in absolute value. Suppose  $\rho \leq \delta/8$ . Let  $\mathcal{I}$  be the set of entries obeying this. For  $i \in \mathcal{I}$  using  $\|\tilde{\boldsymbol{r}}_{\tau}\|_{\ell_{\infty}} \leq 4\rho \leq \delta/4$ , we have

$$|r_{\tau,i}| \ge \delta/2 \implies |r_{\tau,i}| + |r_{\tau,i} - \bar{r}_{\tau,i}| \ge \delta/2 \implies |r_{\tau,i} - \bar{r}_{\tau,i}| \ge \delta/4.$$

Consequently, we find that

$$\|\boldsymbol{r}_{\tau} - \bar{\boldsymbol{r}}_{\tau}\|_{\ell_1} \geq |\mathcal{I}|\delta/4.$$

Converting  $\ell_2$  upper bound on the left hand side to  $\ell_1$ , we obtain

$$c\sqrt{n}\frac{\varepsilon_0\Gamma^3K\sqrt{n\log K}}{\lambda(\boldsymbol{C})}\log(\frac{\Gamma\sqrt{n\log K}}{\rho})\geq |\mathcal{I}|\delta/4.$$

Hence, the total number of errors is at most

$$|\mathcal{I}| \le c' \frac{\varepsilon_0 nK}{\delta} \frac{\Gamma^3 \sqrt{\log K}}{\lambda(C)} \log(\frac{\Gamma \sqrt{n \log K}}{\rho})$$

Third statement – Showing zero error: Pick an input sample x from dataset and its clean version  $\tilde{x}$ . We will argue that  $f(W_{\tau}, x) - f(\tilde{W}_{\tau}, \tilde{x})$  is smaller than  $\delta/4$  when  $\varepsilon_0$  is small enough. We again write

$$|f(\boldsymbol{W}_{\tau}, \boldsymbol{x}) - f(\tilde{\boldsymbol{W}}_{\tau}, \tilde{\boldsymbol{x}})| \le |f(\boldsymbol{W}_{\tau}, \boldsymbol{x}) - f(\tilde{\boldsymbol{W}}_{\tau}, \boldsymbol{x})| + |f(\tilde{\boldsymbol{W}}_{\tau}, \boldsymbol{x}) - f(\tilde{\boldsymbol{W}}_{\tau}, \tilde{\boldsymbol{x}})|$$

The first term can be bounded via

$$|f(\boldsymbol{W}_{\tau}, \boldsymbol{x}) - f(\tilde{\boldsymbol{W}}_{\tau}, \boldsymbol{x})| = |\boldsymbol{v}^{T} \phi(\boldsymbol{W}_{\tau} \boldsymbol{x}) - \boldsymbol{v}^{T} \phi(\tilde{\boldsymbol{W}}_{\tau} \boldsymbol{x})| \le ||\boldsymbol{v}||_{\ell_{2}} ||\phi(\boldsymbol{W}_{\tau} \boldsymbol{x}) - \phi(\tilde{\boldsymbol{W}}_{\tau} \boldsymbol{x})||_{\ell_{2}}$$
(6.79)

$$\leq \Gamma \| \boldsymbol{W}_{\tau} - \tilde{\boldsymbol{W}}_{\tau} \|_{F} \tag{6.80}$$

$$\leq \mathcal{O}(\varepsilon_0 \frac{\Gamma^5 K^2}{\lambda(C)^2} \log(\frac{\Gamma \sqrt{n \log K}}{\rho})^3)$$
(6.81)

Next, we need to bound

$$|f(\tilde{W}_{\tau}, x) - f(\tilde{W}_{\tau}, \tilde{x})| \le |v^T \phi(\tilde{W}_{\tau} x) - v^T \phi(\tilde{W}_{\tau} \tilde{x})|$$
(6.82)

where  $\|\tilde{\boldsymbol{W}}_{\tau} - \boldsymbol{W}_0\|_F \leq \mathcal{O}(\Gamma\sqrt{\frac{K\log K}{\lambda(C)}})$ ,  $\|\boldsymbol{x} - \tilde{\boldsymbol{x}}\|_{\ell_2} \leq \varepsilon_0$  and  $\boldsymbol{W}_0 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \boldsymbol{I})$ . Consequently, using by assumption we have

$$k \ge \mathcal{O}(\|\tilde{\boldsymbol{W}} - \boldsymbol{W}_0\|_F^2) = \mathcal{O}(\Gamma^2 \frac{K \log K}{\lambda(\boldsymbol{C})}),$$

and applying an argument similar to Theorem 2.4 (detailed in Appendix B), with probability at  $1 - 1/n^{100}$ , we find that

$$|f(\tilde{\boldsymbol{W}}_{\tau}, \boldsymbol{x}) - f(\tilde{\boldsymbol{W}}_{\tau}, \tilde{\boldsymbol{x}})| \le C' \Gamma \varepsilon_0 (\|\tilde{\boldsymbol{W}}_{\tau} - \boldsymbol{W}_0\|_F + \sqrt{\log n})$$
(6.83)

$$C\Gamma\varepsilon_0(\Gamma\sqrt{\frac{K\log K}{\lambda(C)}} + \sqrt{\log n}).$$
 (6.84)

Combining the two bounds above we get

$$|f(\boldsymbol{W}_{\tau}, \boldsymbol{x}) - f(\tilde{\boldsymbol{W}}_{\tau}, \tilde{\boldsymbol{x}})| \le \varepsilon_0 \mathcal{O}\left(\frac{\Gamma^5 K^2}{\lambda(\boldsymbol{C})^2} \log\left(\frac{\Gamma\sqrt{n \log K}}{\rho}\right)^3 + \Gamma\left(\Gamma\sqrt{\frac{K \log K}{\lambda(\boldsymbol{C})}} + \sqrt{\log n}\right)\right)$$
(6.85)

$$\leq \varepsilon_0 \mathcal{O}\left(\frac{\Gamma^5 K^2}{\lambda(C)^2} \log\left(\frac{\Gamma \sqrt{n \log K}}{\rho}\right)^3\right).$$
(6.86)

Hence, if  $\varepsilon_0 \leq c' \frac{\delta \lambda(C)^2}{\Gamma^5 K^2 \log(\frac{\Gamma \sqrt{n \log K}}{\rho})^3}$ , we obtain that, for all  $1 \leq i \leq n$ ,

$$|f(\boldsymbol{W}_{\tau},\boldsymbol{x}_{i}) - \tilde{y}_{i}| < |f(\tilde{\boldsymbol{W}}_{\tau},\tilde{\boldsymbol{x}}_{i}) - f(\boldsymbol{W}_{\tau},\boldsymbol{x}_{i})| + |f(\tilde{\boldsymbol{W}}_{\tau},\tilde{\boldsymbol{x}}_{i}) - \tilde{y}_{i}|\tilde{y}_{i}| \le 4\rho + \frac{\delta}{4}.$$

If  $\rho \leq \delta/8$ , we obtain

$$|f(\boldsymbol{W}_{\tau}, \boldsymbol{x}_i) - \tilde{y}_i| < \delta/2$$

hence,  $W_{\tau}$  outputs the correct decision for all samples.

Fourth statement - Distance: This follows from the triangle inequality

$$\| \boldsymbol{W}_{\tau} - \boldsymbol{W}_{0} \|_{F} \le \| \boldsymbol{W}_{\tau} - \tilde{\boldsymbol{W}}_{\tau} \|_{F} + \| \tilde{\boldsymbol{W}}_{\tau} - \boldsymbol{W}_{0} \|_{F}$$

We have that right hand side terms are at most  $\mathcal{O}(\Gamma\sqrt{\frac{K\log K}{\lambda(C)}})$  and  $\mathcal{O}(t\eta\varepsilon_0\frac{\Gamma^4Kn}{\lambda(C)}\log(\frac{\Gamma\sqrt{n\log K}}{\rho})^2)$  from Theorems 6.12 and 2.3 respectively. This implies (6.75).

## Acknowledgements

M. Soltanolkotabi is supported by the Packard Fellowship in Science and Engineering, a Sloan Research Fellowship in Mathematics, an NSF-CAREER under award #1846369, the Air Force Office of Scientific Research Young Investigator Program (AFOSR-YIP) under award #FA9550-18-1-0078, an NSF-CIF award #1813877, and a Google faculty research award.

## References

- [1] ABBE, E., BANDEIRA, A. S., AND HALL, G. Exact recovery in the stochastic block model. *IEEE Transactions on Information Theory* 62, 1 (2016), 471–487.
- [2] Allen-Zhu, Z., Li, Y., and Liang, Y. Learning and generalization in overparameterized neural networks, going beyond two layers. arXiv preprint arXiv:1811.04918 (2018).
- [3] Allen-Zhu, Z., Li, Y., and Song, Z. A convergence theory for deep learning via over-parameterization. arXiv preprint arXiv:1811.03962 (2018).
- [4] Arora, S., Cohen, N., and Hazan, E. On the optimization of deep networks: Implicit acceleration by overparameterization. arXiv preprint arXiv:1802.06509 (2018).
- [5] Arora, S., Ge, R., Neyshabur, B., and Zhang, Y. Stronger generalization bounds for deep nets via a compression approach.
- [6] BALAKRISHNAN, S., Du, S. S., Li, J., and Singh, A. Computationally efficient robust sparse estimation in high dimensions. In *Conference on Learning Theory* (2017), pp. 169–212.
- [7] BARTLETT, P., FOSTER, D. J., AND TELGARSKY, M. Spectrally-normalized margin bounds for neural networks.
- [8] Belkin, M., Hsu, D., and Mitra, P. Overfitting or perfect fitting? risk bounds for classification and regression rules that interpolate.
- [9] Belkin, M., Rakhlin, A., and Tsybakov, A. B. Does data interpolation contradict statistical optimality?
- [10] Bhatia, K., Jain, P., and Kar, P. Robust regression via hard thresholding. In *Advances in Neural Information Processing Systems* (2015), pp. 721–729.
- [11] Brutzkus, A., and Globerson, A. Over-parameterization improves generalization in the xor detection problem.
- [12] Brutzkus, A., Globerson, A., Malach, E., and Shalev-Shwartz, S. Sgd learns over-parameterized networks that provably generalize on linearly separable data. arXiv preprint arXiv:1710.10174 (2017).
- [13] Brutzkus, A., Globerson, A., Malach, E., and Shalev-Shwartz, S. Sgd learns over-parameterized networks that provably generalize on linearly separable data.
- [14] CAI, T. T., LI, X., ET AL. Robust and computationally feasible community detection in the presence of arbitrary outlier nodes. *The Annals of Statistics* 43, 3 (2015), 1027–1059.
- [15] CANDÈS, E. J., LI, X., MA, Y., AND WRIGHT, J. Robust principal component analysis? *Journal of the ACM (JACM) 58*, 3 (2011), 11.
- [16] CHAUDHARI, P., CHOROMANSKA, A., SOATTO, S., LECUN, Y., BALDASSI, C., BORGS, C., CHAYES, J., SAGUN, L., AND ZECCHINA, R. Entropy-sgd: Biasing gradient descent into wide valleys. arXiv preprint arXiv:1611.01838 (2016).
- [17] CHEN, Y., CARAMANIS, C., AND MANNOR, S. Robust sparse regression under adversarial corruption. In *International Conference on Machine Learning* (2013), pp. 774–782.
- [18] CHIZAT, L., AND BACH, F. On the global convergence of gradient descent for over-parameterized models using optimal transport. arXiv preprint arXiv:1805.09545 (2018).

- [19] Diakonikolas, I., Kamath, G., Kane, D. M., Li, J., Steinhardt, J., and Stewart, A. Sever: A robust meta-algorithm for stochastic optimization. arXiv preprint arXiv:1803.02815 (2018).
- [20] Du, S. S., Lee, J. D., Li, H., Wang, L., and Zhai, X. Gradient descent finds global minima of deep neural networks. arXiv preprint arXiv:1811.03804 (2018).
- [21] Du, S. S., Zhai, X., Poczos, B., and Singh, A. Gradient descent provably optimizes over-parameterized neural networks. arXiv preprint arXiv:1810.02054 (2018).
- [22] FOYGEL, R., AND MACKEY, L. Corrupted sensing: Novel guarantees for separating structured signals. *IEEE Transactions on Information Theory* 60, 2 (2014), 1223–1247.
- [23] FRÉNAY, B., KABÁN, A., ET AL. A comprehensive introduction to label noise. In ESANN (2014).
- [24] GOLOWICH, N., RAKHLIN, A., AND SHAMIR, O. Size-independent sample complexity of neural networks.
- [25] Gur-Ari, G., Roberts, D. A., and Dyer, E. Gradient descent happens in a tiny subspace. arXiv preprint arXiv:1812.04754 (2018).
- [26] HAN, B., YAO, Q., YU, X., NIU, G., XU, M., HU, W., TSANG, I., AND SUGIYAMA, M. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In *Advances in Neural Information Processing Systems* (2018), pp. 8536–8546.
- [27] Hochreiter, S., and Schmidhuber, J. Flat minima. Neural Computation 9, 1 (1997), 1–42.
- [28] JI, Z., AND TELGARSKY, M. Gradient descent aligns the layers of deep linear networks.
- [29] KESKAR, N. S., MUDIGERE, D., NOCEDAL, J., SMELYANSKIY, M., AND TANG, P. T. P. On large-batch training for deep learning: Generalization gap and sharp minima. arXiv preprint arXiv:1609.04836 (2016).
- [30] KHETAN, A., LIPTON, Z. C., AND ANANDKUMAR, A. Learning from noisy singly-labeled data. arXiv preprint arXiv:1712.04577 (2017).
- [31] KLIVANS, A., KOTHARI, P. K., AND MEKA, R. Efficient algorithms for outlier-robust regression. arXiv preprint arXiv:1803.03241 (2018).
- [32] Li, X. Compressed sensing and matrix completion with constant proportion of corruptions. *Constructive Approximation* 37, 1 (2013), 73–99.
- [33] Li, Y., and Liang, Y. Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in Neural Information Processing Systems* (2018), pp. 8168–8177.
- [34] LIANG, T., AND RAKHLIN, A. Just interpolate: Kernel "ridgeless" regression can generalize.
- [35] Liu, L., Shen, Y., Li, T., and Caramanis, C. High dimensional robust sparse regression. arXiv preprint arXiv:1805.11643 (2018).
- [36] Malach, E., and Shalev-Shwartz, S. Decoupling" when to update from how to update". In Advances in Neural Information Processing Systems (2017), pp. 960–970.
- [37] MENON, A. K., VAN ROOYEN, B., AND NATARAJAN, N. Learning from binary labels with instance-dependent noise. *Machine Learning* (2018), 1–35.
- [38] OYMAK, S., AND SOLTANOLKOTABI, M. Overparameterized nonlinear learning: Gradient descent takes the shortest path? arXiv preprint arXiv:1812.10004 (2018).
- [39] OYMAK, S., AND SOLTANOLKOTABI, M. Overparameterized nonlinear learning: Gradient descent takes the shortest path? arXiv preprint arXiv:1812.10004 (2018).

- [40] OYMAK, S., AND SOLTANOLKOTABI, M. Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. arXiv preprint arXiv:1902.04674 (2019).
- [41] Prasad, A., Suggala, A. S., Balakrishnan, S., and Ravikumar, P. Robust estimation via robust gradient estimation. arXiv preprint arXiv:1802.06485 (2018).
- [42] REED, S., LEE, H., ANGUELOV, D., SZEGEDY, C., ERHAN, D., AND RABINOVICH, A. Training deep neural networks on noisy labels with bootstrapping. arXiv preprint arXiv:1412.6596 (2014).
- [43] Ren, M., Zeng, W., Yang, B., and Urtasun, R. Learning to reweight examples for robust deep learning. arXiv preprint arXiv:1803.09050 (2018).
- [44] ROLNICK, D., VEIT, A., BELONGIE, S., AND SHAVIT, N. Deep learning is robust to massive label noise. arXiv preprint arXiv:1705.10694 (2017).
- [45] SAGUN, L., EVCI, U., GUNEY, V. U., DAUPHIN, Y., AND BOTTOU, L. Empirical analysis of the hessian of over-parametrized neural networks. arXiv preprint arXiv:1706.04454 (2017).
- [46] Schur, J. Bemerkungen zur theorie der beschränkten bilinearformen mit unendlich vielen veränderlichen. Journal für die reine und angewandte Mathematik 140 (1911), 1–28.
- [47] Scott, C., Blanchard, G., and Handy, G. Classification with asymmetric label noise: Consistency and maximal denoising. In *Conference On Learning Theory* (2013), pp. 489–511.
- [48] Shen, Y., and Sanghavi, S. Iteratively learning from the best. arXiv preprint arXiv:1810.11874 (2018).
- [49] SOLTANOLKOTABI, M., JAVANMARD, A., AND LEE, J. D. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information Theory* (2018).
- [50] Song, M., Montanari, A., and Nguyen, P. A mean field view of the landscape of two-layers neural networks. In *Proceedings of the National Academy of Sciences* (2018), vol. 115, pp. E7665–E7671.
- [51] SOUDRY, D., AND CARMON, Y. No bad local minima: Data independent training error guarantees for multilayer neural networks.
- [52] SOUDRY, D., HOFFER, E., NACSON, M. S., GUNASEKAR, S., AND SREBRO, N. The implicit bias of gradient descent on separable data. arXiv preprint arXiv:1710.10345 (2017).
- [53] Venturi, L., Bandeira, A., and Bruna, J. Spurious valleys in two-layer neural network optimization landscapes. arXiv preprint arXiv:1802.06384 (2018).
- [54] VINAYAK, R. K., OYMAK, S., AND HASSIBI, B. Graph clustering with missing data: Convex algorithms and analysis. In *Advances in Neural Information Processing Systems* (2014), pp. 2996–3004.
- [55] XIE, B., LIANG, Y., AND SONG, L. Diverse neural network learns true target functions. arXiv preprint arXiv:1611.03131 (2016).
- [56] Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning requires rethinking generalization. *International Conference on Learning Representations* (2016).
- [57] ZHANG, Z., AND SABUNCU, M. R. Generalized cross entropy loss for training deep neural networks with noisy labels. arXiv preprint arXiv:1805.07836 (2018).
- [58] Zhu, Z., Soudry, D., Eldar, Y. C., and Wakin, M. B. The global optimization geometry of shallow linear neural networks.
- [59] ZOU, D., CAO, Y., ZHOU, D., AND GU, Q. Stochastic gradient descent optimizes over-parameterized deep relu networks. arXiv preprint arXiv:1811.08888 (2018).

## A Proof of Lemma 2.5

Create two matrices  $X \in \mathbb{R}^{s \times d}$  and  $\tilde{X} \in \mathbb{R}^{s \times d}$  by concatenating the input samples. Note that the matrix  $X - \tilde{X}$  has i.i.d.  $\mathcal{N}(0, 2\varepsilon_0^2/d)$  entries. Thus, using standard results regarding the concentration of the spectral norm with probability at least  $1 - e^{-d/2}$ , we have

$$\|\boldsymbol{X} - \tilde{\boldsymbol{X}}\| \le \sqrt{2} \left(\sqrt{\frac{s}{d}} + 2\right) \varepsilon_0 \le 5\varepsilon_0.$$

Define the vectors  $\boldsymbol{y}, \tilde{\boldsymbol{y}} \in \mathbb{R}^s$  with entries given by  $y_i$  and  $\tilde{y}_i$ , respectively. Suppose  $\boldsymbol{W}$  fits these labels perfectly. Using the fact that  $\|\boldsymbol{v}\|_{\ell_2} = 1$ , we can conclude that

$$\sqrt{s}\delta \leq \|\boldsymbol{y} - \tilde{\boldsymbol{y}}\|_{\ell_{2}} = \|f(\boldsymbol{W}, \boldsymbol{X}) - f(\boldsymbol{W}, \tilde{\boldsymbol{X}})\|_{\ell_{2}},$$

$$= \|\boldsymbol{v}^{T}(\phi(\boldsymbol{W}\boldsymbol{X}) - \phi(\boldsymbol{W}\tilde{\boldsymbol{X}}))\|_{\ell_{2}},$$

$$\leq \Gamma \|\boldsymbol{v}\|_{\ell_{2}} \|\boldsymbol{W}(\boldsymbol{X} - \tilde{\boldsymbol{X}})\|_{F},$$

$$\leq \Gamma \|\boldsymbol{X} - \tilde{\boldsymbol{X}}\| \|\boldsymbol{W}\|_{F} \leq 5\Gamma \varepsilon_{0} \|\boldsymbol{W}\|_{F}.$$

This implies the desired lower bound on  $\|\mathbf{W}\|_F$ .

## B Single label perturbation

Note that

$$|f(\boldsymbol{W}, \boldsymbol{x}) - f(\boldsymbol{W}, \widetilde{\boldsymbol{x}})| = |\boldsymbol{v}^{T} (\phi (\boldsymbol{W} \boldsymbol{x}) - \phi (\boldsymbol{W} \widetilde{\boldsymbol{x}}))|$$

$$\leq |\boldsymbol{v}^{T} (\phi (\boldsymbol{W} \boldsymbol{x}) - \phi (\boldsymbol{W} \widetilde{\boldsymbol{x}})) - \boldsymbol{v}^{T} (\phi (\boldsymbol{W}_{0} \boldsymbol{x}) - \phi (\boldsymbol{W}_{0} \widetilde{\boldsymbol{x}}))| + |\boldsymbol{v}^{T} (\phi (\boldsymbol{W}_{0} \boldsymbol{x}) - \phi (\boldsymbol{W}_{0} \widetilde{\boldsymbol{x}}))|$$

To continue note that by the general mean value theorem there exists a point  $(\overline{W}, \overline{x})$  in the square  $(W_0, x), (W_0, \widetilde{x}), (W, x)$ , and  $(W, \widetilde{x})$  such that

$$(f(\boldsymbol{W}, \boldsymbol{x}) - f(\boldsymbol{W}_0, \boldsymbol{x})) - (f(\boldsymbol{W}, \widetilde{\boldsymbol{x}}) - f(\boldsymbol{W}_0, \widetilde{\boldsymbol{x}}))$$

$$= (\boldsymbol{x} - \widetilde{\boldsymbol{x}})^T \overline{\boldsymbol{W}}^T \operatorname{diag}(\boldsymbol{v}) \operatorname{diag}(\phi''(\overline{\boldsymbol{W}}\overline{\boldsymbol{x}})) (\boldsymbol{W} - \boldsymbol{W}_0) \overline{\boldsymbol{x}} + \boldsymbol{v}^T \operatorname{diag}(\phi'(\overline{\boldsymbol{W}}\overline{\boldsymbol{x}})) (\boldsymbol{W} - \boldsymbol{W}_0) (\boldsymbol{x} - \widetilde{\boldsymbol{x}})$$

Using the above we have that

$$\left| (f(\boldsymbol{W}, \boldsymbol{x}) - f(\boldsymbol{W}_{0}, \boldsymbol{x})) - (f(\boldsymbol{W}, \widetilde{\boldsymbol{x}}) - f(\boldsymbol{W}_{0}, \widetilde{\boldsymbol{x}})) \right| \stackrel{(a)}{\leq} \left| (\boldsymbol{x} - \widetilde{\boldsymbol{x}})^{T} \overline{\boldsymbol{W}}^{T} \operatorname{diag}(\boldsymbol{v}) \operatorname{diag}(\boldsymbol{\phi}''(\overline{\boldsymbol{W}} \overline{\boldsymbol{x}})) (\boldsymbol{W} - \boldsymbol{W}_{0}) \boldsymbol{x} \right| \\
+ \left| \boldsymbol{v}^{T} \operatorname{diag}(\boldsymbol{\phi}'(\overline{\boldsymbol{W}} \overline{\boldsymbol{x}})) (\boldsymbol{W} - \boldsymbol{W}_{0}) (\boldsymbol{x} - \widetilde{\boldsymbol{x}}) \right| \\
\stackrel{(b)}{\leq} \left( \|\boldsymbol{v}\|_{\ell_{\infty}} \|\overline{\boldsymbol{x}}\|_{\ell_{2}} \|\overline{\boldsymbol{W}}\| + \|\boldsymbol{v}\|_{\ell_{2}} \right) \Gamma \|\boldsymbol{x} - \widetilde{\boldsymbol{x}}\|_{\ell_{2}} \|\boldsymbol{W} - \boldsymbol{W}_{0}\| \\
\stackrel{(c)}{\leq} \left( \frac{1}{\sqrt{k}} \|\overline{\boldsymbol{W}}\| + 1 \right) \Gamma \|\boldsymbol{x} - \widetilde{\boldsymbol{x}}\|_{\ell_{2}} \|\boldsymbol{W} - \boldsymbol{W}_{0}\| \\
\stackrel{(e)}{\leq} \left( \frac{1}{\sqrt{k}} \|\boldsymbol{W}_{0}\| + \frac{1}{\sqrt{k}} \|\overline{\boldsymbol{W}} - \boldsymbol{W}_{0}\| + 1 \right) \Gamma \|\boldsymbol{x} - \widetilde{\boldsymbol{x}}\|_{\ell_{2}} \|\boldsymbol{W} - \boldsymbol{W}_{0}\| \\
\stackrel{(f)}{\leq} \left( \frac{1}{\sqrt{k}} \|\boldsymbol{W}_{0}\| + \frac{1}{\sqrt{k}} \|\overline{\boldsymbol{W}} - \boldsymbol{W}_{0}\| + 1 \right) \Gamma \|\boldsymbol{x} - \widetilde{\boldsymbol{x}}\|_{\ell_{2}} \|\boldsymbol{W} - \boldsymbol{W}_{0}\| \\
\stackrel{(g)}{\leq} \left( \frac{1}{\sqrt{k}} \|\overline{\boldsymbol{W}} - \boldsymbol{W}_{0}\|_{F} + 3 + 2\sqrt{\frac{d}{k}} \right) \Gamma \|\boldsymbol{x} - \widetilde{\boldsymbol{x}}\|_{\ell_{2}} \|\boldsymbol{W} - \boldsymbol{W}_{0}\| \\
\stackrel{(g)}{\leq} C\Gamma \|\boldsymbol{x} - \widetilde{\boldsymbol{x}}\|_{\ell_{2}} \|\boldsymbol{W} - \boldsymbol{W}_{0}\| \tag{B.1}$$

Here, (a) follows from the triangle inequality, (b) from simple algebraic manipulations along with the fact that  $|\phi'(z)| \leq \Gamma$  and  $|\phi''(z)| \leq \Gamma$ , (c) from the fact that  $\mathbf{v}_{\ell} = \pm \frac{1}{\sqrt{k}}$ , (d) from  $\|\mathbf{x}\|_{\ell_2} = \|\widetilde{\mathbf{x}}\|_{\ell_2} = 1$  which implies  $\|\overline{\mathbf{x}}\|_{\ell_2} \leq 1$ , (e) from triangular inequality, (f) from the fact that Frobenius norm dominates the spectral norm, (g) from the fact that with probability at least  $1 - 2e^{-(d+k)}$ ,  $\|\mathbf{W}_0\| \leq 2(\sqrt{k} + \sqrt{d})$ , and (h) from the fact that  $\|\overline{\mathbf{W}} - \mathbf{W}_0\| \leq \|\mathbf{W} - \mathbf{W}_0\|_F \leq \widetilde{c}\sqrt{k}$  and  $k \geq cd$ .

Next we note that for a Gaussian random vector  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  we have

$$\|\phi(\boldsymbol{g}^{T}\boldsymbol{x}) - \phi(\boldsymbol{g}^{T}\widetilde{\boldsymbol{x}})\|_{\psi_{2}} = \|\phi(\boldsymbol{g}^{T}\boldsymbol{x}) - \phi(\boldsymbol{g}^{T}\widetilde{\boldsymbol{x}})\|_{\psi_{2}}$$

$$= \|\phi'\left(t\boldsymbol{g}^{T}\boldsymbol{x} + (1-t)\boldsymbol{g}^{T}\widetilde{\boldsymbol{x}}\right)\boldsymbol{g}^{T}(\boldsymbol{x} - \widetilde{\boldsymbol{x}})\|_{\psi_{2}}$$

$$\leq \Gamma \|\boldsymbol{g}^{T}(\boldsymbol{x} - \widetilde{\boldsymbol{x}})\|_{\psi_{2}}$$

$$\leq c\Gamma \|\boldsymbol{x} - \widetilde{\boldsymbol{x}}\|_{\ell_{2}}.$$
(B.2)

Also note that

$$f(\boldsymbol{W}_{0}, \boldsymbol{x}) - f(\boldsymbol{W}_{0}, \widetilde{\boldsymbol{x}}) = \boldsymbol{v}^{T} \left( \phi \left( \boldsymbol{W}_{0} \boldsymbol{x} \right) - \phi \left( \boldsymbol{W}_{0} \widetilde{\boldsymbol{x}} \right) \right)$$
$$\sim \sum_{\ell=1}^{k} \boldsymbol{v}_{\ell} \left( \phi \left( \boldsymbol{g}_{\ell}^{T} \boldsymbol{x} \right) - \phi \left( \boldsymbol{g}_{\ell}^{T} \widetilde{\boldsymbol{x}} \right) \right)$$

where  $g_1, g_2, ..., g_k$  are i.i.d. vectors with  $\mathcal{N}(0, I_d)$  distribution. Also for  $\boldsymbol{v}$  obeying  $\mathbf{1}^T \boldsymbol{v} = 0$  this random variable has mean zero. Hence, using the fact that weighted sum of subGaussian random variables are subgaussian combined with  $(\mathbf{B}.2)$  we conclude that  $f(\boldsymbol{W}_0, \boldsymbol{x}) - f(\boldsymbol{W}_0, \widetilde{\boldsymbol{x}})$  is also subGaussian obeying  $\|f(\boldsymbol{W}_0, \boldsymbol{x}) - f(\boldsymbol{W}_0, \widetilde{\boldsymbol{x}})\|_{\psi_2} \le c\Gamma \|\boldsymbol{v}\|_{\ell_2} \|\boldsymbol{x} - \widetilde{\boldsymbol{x}}\|_{\ell_2}$ . Thus

$$|f(\boldsymbol{W}_{0},\boldsymbol{x}) - f(\boldsymbol{W}_{0},\widetilde{\boldsymbol{x}})| \le ct\Gamma \|\boldsymbol{v}\|_{\ell_{2}} \|\boldsymbol{x} - \widetilde{\boldsymbol{x}}\|_{\ell_{2}} = ct\Gamma \|\boldsymbol{x} - \widetilde{\boldsymbol{x}}\|_{\ell_{2}},$$
(B.3)

with probability at least  $1 - e^{-\frac{t^2}{2}}$ . Thus, using  $t = 2\sqrt{\log n}$  for n data points

$$|f(\boldsymbol{W}_0, \boldsymbol{x}_i) - f(\boldsymbol{W}_0, \widetilde{\boldsymbol{x}}_i)| \le 2c\Gamma\sqrt{\log n} \|\boldsymbol{x}_i - \widetilde{\boldsymbol{x}}_i\|_{\ell_2},$$

holds for all i = 1, 2, ..., n with probability at least

$$1 - ne^{-\frac{t^2}{2}} \ge 1 - \frac{1}{n^{100}}.$$