

---

# Supplementary material: what training reveals about neural network complexity

---

**Andreas Loukas**  
EPFL  
andreas.loukas@epfl.ch

**Marinos Poiitis**  
Aristotle University of Thessaloniki  
mpoiitis@csd.auth.gr

**Stefanie Jegelka**  
MIT  
stefje@mit.edu

The supplementary material commences in Section B by presenting the proofs supporting our theoretical claims. Section A discusses experimental details and shows additional results. Finally, Section C lays out theoretical results of side-interest.

## A Additional empirical results

### A.1 Description of Task 1

The input data of Task 1 are generated by the following two step procedure:

First, we sample  $N = 100$  points  $\mathbf{z}_i \in [-1, 1]^2$  uniformly at random and assign them a ground truth label according to the sinusoidal function:

$$y_i = \cos(2\pi\omega \mathbf{z}_i(1)) \cdot \cos(2\pi\omega \mathbf{z}_i(2)) \in [-1, 1],$$

where  $\omega$  is interpreted as a frequency and we set  $\omega \in \{0.25, 0.5, 0.75, 1.0\}$  in our experiments. The four resulting functions are visualized in Figure A.2.

We then determine  $\{\mathbf{x}_i\}_{i=1}^N$  by isometrically embedding  $\{\mathbf{z}_i\}_{i=1}^N$  into  $\mathbb{R}^{10}$ . We achieve this by selecting the first 2 columns  $\mathbf{R} \in \mathbb{R}^{10 \times 2}$  of a random  $10 \times 10$  unitary matrix and setting  $\mathbf{x}_i = \mathbf{R} \mathbf{z}_i$ . This procedure ensures that the distances between points remains the same in high dimensions.

### A.2 Distance to initialization for Task 2

We focus on the image classification CNN trained with a BCE loss. Figure A.1 depicts the distance from initialization  $\|\mathbf{b}_1^{(t)} - \mathbf{b}_1^{(0)}\|_2$  in the last 10 training epochs.

As explained in Section 4.2, when a BCE loss is utilised, the derivative of the loss becomes unbounded which stops Corollary 3 from applying. Interestingly, Figure A.1 confirms this by showing that the distance is not an increasing function of complexity. The reverse phenomenon can be observed when an MSE loss is utilized (see Figure 3b).

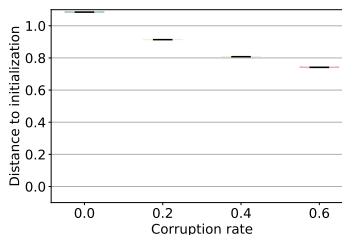


Figure A.1: Distance to init. with BCE loss.

### A.3 Visualizing linear regions

Aiming to gain intuition about the behavior of NNs in linear regions close and far the training data, we take a closer look at the function an MLP is learning when trained to solve task 1 ( $\omega = 0.5$ , 2 hidden layers,  $N = 200$ ).

Figures A.3a and A.3b depict, respectively, the real and learned function projected in 2D (recall that the true function is isometrically embedded in 10D). Blue dots are training data points. The

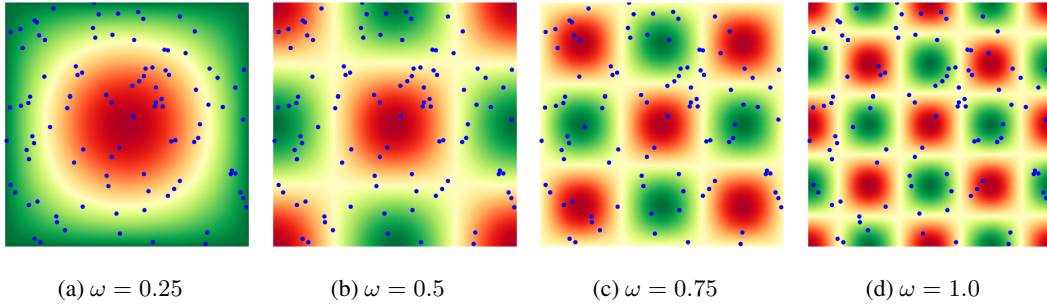


Figure A.2: The surface of the sinusoidal function from where the input points are sampled, for different frequencies  $\omega$ . Sampled points are plotted on top of the surface.

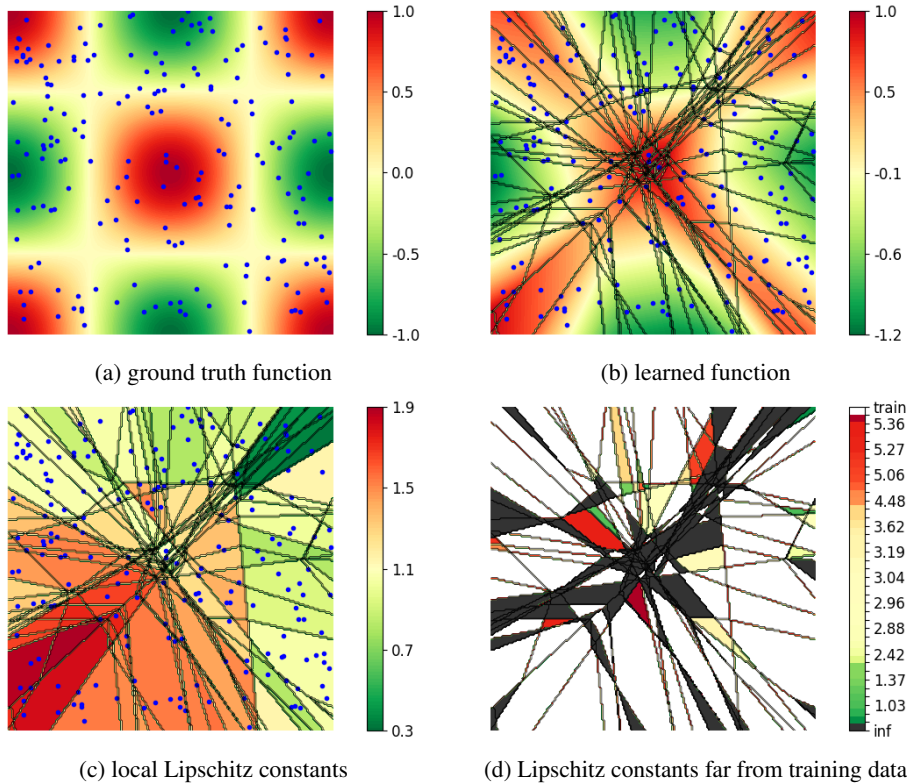


Figure A.3: Visual illustration of the linear regions (top right) of the trained NN when fitting a sinusoidal function (top left), together with their local Lipschitz constants (bottom left) and those far from the training data as predicted by Theorem 2 (bottom right).

boundaries between region are indicated with black lines. As observed, there is a large number of linear regions of varying sizes with the smaller and more densely packed regions being found close to the  $(0,0)$  point.

The bottom two panels display the local Lipschitz constants (i.e., the magnitude of the gradient within every region). In Figure A.3c we can see the real constants at all regions. Interestingly, it appears that low- and high-Lipschitz constants are clustered, which likely follows from the hierarchical region formation process: in other words, regions within the same cluster fall within the same region of a shallower sub-network and are split by a higher layer.

Figure A.3d distinguishes between regions containing training points (in white) and the rest (in color). We color empty regions depending on the bound given by Theorem 2 and black regions are those for which the theorem does not have predictive power. We observe that, though the proposed theory

allows us to make statements about the function behavior far from the training data, the theory does not explain the global behavior of the NN. This motivates the introduction of Dropout in the analysis of Section 5.2: by exploiting stochasticity we can infer more properties about the NN complexity from the training trajectory. Intuitively, using Dropout during a sufficiently long training, one is able to deduce from the observed bias updates (specifically vector  $\varphi_T$  in Theorem 2) the Lipschitz constants within more regions (thus they would also bound the Lipschitz constant within some of the non-white regions in Fig A.3d). Moreover, though encountering each and every region in the training would likely take a very long time, Theorem 2 implies that only a small subset of regions suffice to approximate the global Lipschitz constant up to a logarithmic factor.

We finally observe that, in the regions where it applies, Theorem 2 yields a bound that is a constant factor away from the real local Lipschitz constants: the bound overestimates the constants by roughly a factor of four.

#### A.4 Total trajectory length

Figure A.4 displays the length of the entire normalized bias trajectory at every point in the training. Thus, Figure A.4 corresponds to the integral of Figures 2b and 2e, which focus on the length of the normalized bias trajectory within every epoch. We note also that all NNs have been trained until they could closely fit the training set.

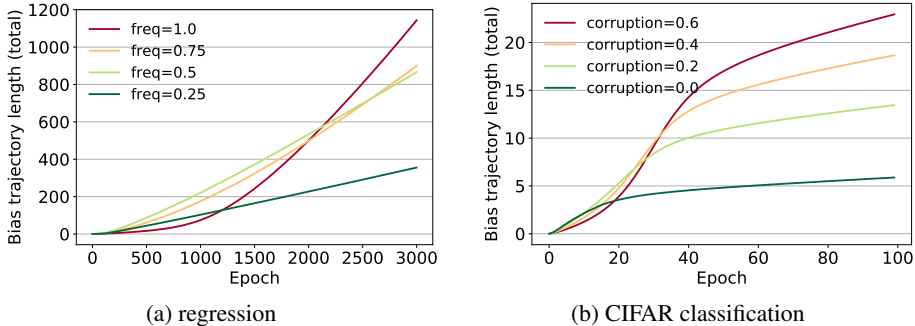


Figure A.4: Length of the normalized trajectory at every pointing in the training starting from initialization. As expected, the trajectory length of NNs grows with the complexity of the function they are learning.

A side-by-side comparison with Figures 2a and 2d reveals that, between any two NNs that have fitted the training data equally well, the one that implements a higher complexity function has consistently a longer trajectory.

#### A.5 Effect of architecture on bias trajectory

We next evaluate the effect of the NN architecture on the optimization trajectory. We focus on the MNIST dataset [1] and train an MLP and a CNN to distinguish between digits ‘3’ and ‘6’ based on a training set consisting of 100 and 1000 images per class. For consistency with the previous experiments, we used the same NN architectures for the MLP and CNN as those employed for Tasks 1 and 2, respectively (though both NNs now feature a sigmoid activation in the last layer). The networks are trained using SGD with a BCE loss and a learning rate of  $\alpha_t = 0.002$ .

Figure A.5 depicts the training loss, normalized bias trajectory length, and test loss for each dataset. Note that, in contrast to Figure 2, here all three measures are computed over time-intervals of 100 iterations (rather than per epoch). As expected, when the training set is small, both architectures fit the training data equally well after roughly 20k iterations, but the CNN overfits less. By observing the length of the bias trajectory, we deduce that the MLP is learning a more complex function than the CNN. Thus, in the MNIST100 case, there is a correlation between trajectory length and generalization with the NN architecture that is more appropriate for the task exhibiting a shorter trajectory.

It is important to remark that the complexity of the learned function is not the sole factor driving generalization (though it can be a crucial factor all other things being equal). Convolutional layers

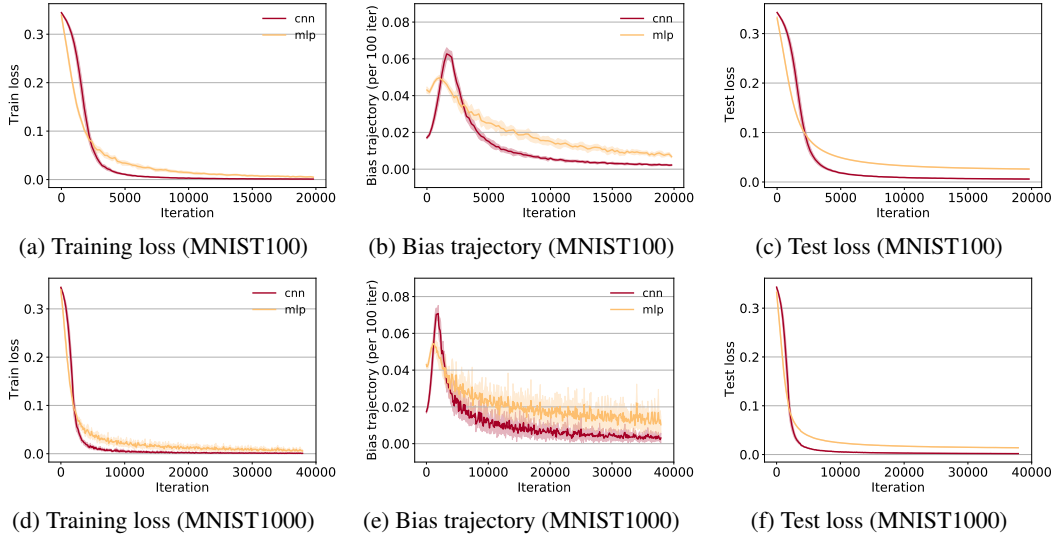


Figure A.5: Training and test behavior of a MLP and a CNN solving a binary MNIST image classification task with different number of samples (100 and 1000 images per class). All three measures (training loss, bias trajectory length, and test loss) are plotted at intervals of 100 iterations.

are indeed more constrained and better suited to image data than fully convolutional ones – thus it is reasonable to expect better generalization than MLPs. Nevertheless, our experiment shows that the CNN, beyond having the right architecture for the task, also learns a slightly lower complexity function than the MLP while fitting the training data equally well or better. Thus, here we mainly use the bias trajectory length as a diagnostic tool that helps us understand what functions the two architectures are learning.

## A.6 Effect of batch size on bias trajectory

This experiment investigates the effect of different batch sizes on the bias trajectory. We adopt the same setup as that of Task 1 (specifically  $\omega = 0.25$  and  $0.75$  in Figure A.6) and train NNs with SGD using batch sizes of 16 and 32, whereas our original experiment used a batch size of 1.

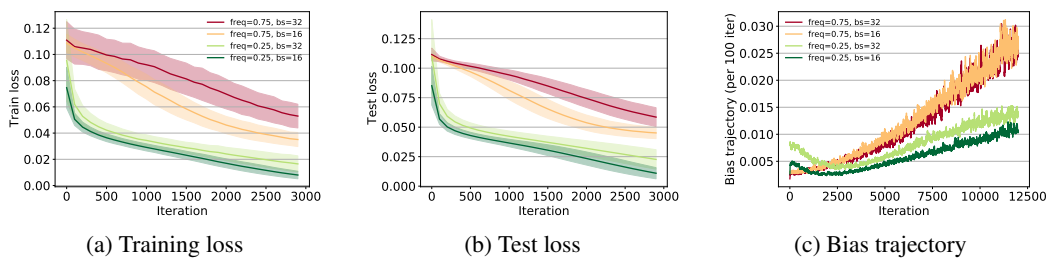


Figure A.6: Training and test behavior of MLPs using different batch size (16 and 32) and evaluated on data sampled from two different frequencies ( $\omega = 0.25$  and  $0.75$  in Task 1). All three measures (training loss, bias trajectory length, and test loss) are plotted at intervals of 100 iterations. The trajectory length correlates with the NN’s complexity for different batch sizes.

The results are consistent with those of Figure 2, with a longer bias trajectory indicating that the NN is fitting a more complex hypothesis and correlating with higher test loss. Increasing the batch size from 1 to 32 also leads to a slight increase in trajectory length, though we currently lack mathematical evidence that support this empirical observation.

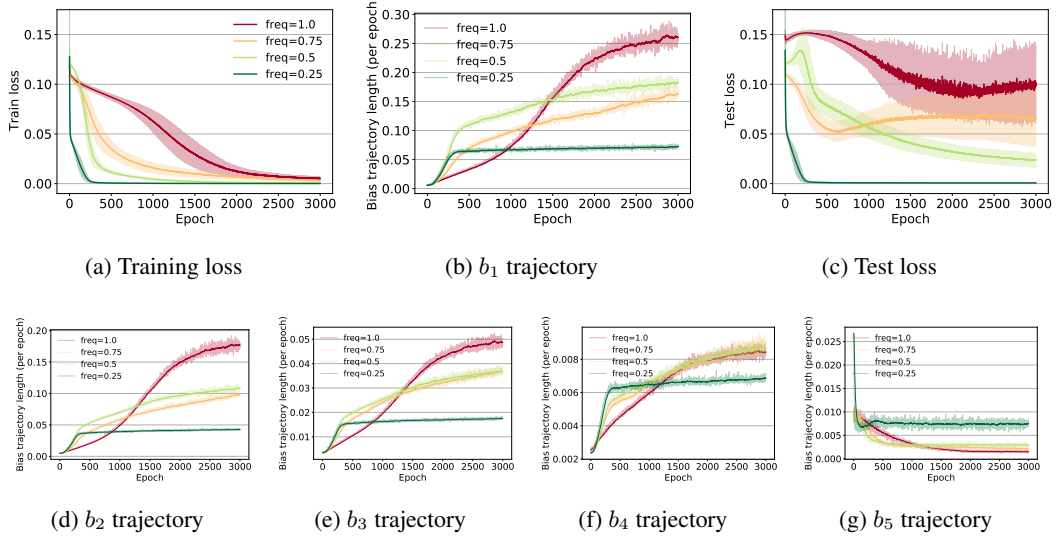


Figure A.7: Illustration of how the biases of all layers change when training a MLP to fit a function of increasing spatial frequency. Sub-figures (a) and (c) show the behavior of the training and test loss, whereas sub-figures (b) and (d-g) depict the per epoch bias trajectory.

## A.7 The trajectory of higher layer biases

Our last experiment examines the training dynamics associated with the biased of higher layers. We focus on Task 1 and replicate the experiment described in Section 6, but now we track the normalized bias trajectory length of all biases and optimize  $\mathbf{W}_1$  freely.

Figure A.7 reports the obtained results. It can be observed that the bias dynamics correlate with task complexity for the first three layers, while being uncorrelated for the last two. To interpret these results, we recall that in the proof of Lemma 1 the trajectory length of  $\mathbf{b}_l$  relates to the Lipschitz constant of subnetwork  $f_{d \leftarrow l+1} = f_d \circ \dots \circ f_{l+1}$  close to the training data. Then, noticing how the trajectory length decays by almost an order of magnitude at each layer, we may infer that the Lipschitz constant of  $f_{d \leftarrow l+1}$  decays quickly with  $l$ . The latter implies that the NN predominantly employs the first few layers to solve the task, whereas the last two layers implement very simple functions.

## B Deferred technical arguments

### B.1 A simple Lemma

**Lemma 1.** Let  $f^{(t)}$  be a  $d$ -layer NN at the  $t$ -th SGD iteration, denote by  $\mathbf{x}^{(t)} \in X$  the point of the training set sampled at that iteration, and set

$$\epsilon_{f^{(t)}}(\mathbf{x}, y) := \left| \frac{\partial \ell(\hat{y}, y)}{\partial \hat{y}} \right|_{\hat{y}=f^{(t)}(\mathbf{x})}. \quad (1)$$

The Lipschitz constant of  $f^{(t)}$  at  $\mathcal{R}_{\mathbf{x}^{(t)}}$  is

$$\frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\alpha_t \cdot \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \cdot \sigma_n(\mathbf{W}_1^{(t)}) \leq \lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}^{(t)}}) \leq \frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\alpha_t \cdot \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \cdot \sigma_1(\mathbf{W}_1^{(t)}),$$

where  $\sigma_1(\mathbf{W}_1^{(t)}) \geq \dots \geq \sigma_n(\mathbf{W}_1^{(t)}) > 0$  are the singular values of  $\mathbf{W}_1^{(t)}$ .

Let us start with some basics. By the chain rule, we have

$$\frac{\partial \ell(f(\mathbf{x}^{(t)}, \mathbf{w}^{(t)}), y^{(t)})}{\partial \mathbf{w}^{(t)}} = \frac{\partial \ell(\hat{y}, y^{(t)})}{\partial \hat{y}} \cdot \frac{\partial f(\mathbf{x}^{(t)}, \mathbf{w}^{(t)})}{\partial \mathbf{w}^{(t)}}$$

with  $\hat{y} = f(\mathbf{x}^{(t)}, \mathbf{w}^{(t)})$ , whereas the gradient w.r.t. the bias of the  $\ell$ -th layer is given by

$$\left( \frac{\partial f(\mathbf{x}^{(t)}, \mathbf{w}^{(t)})}{\partial \mathbf{b}_l^{(t)}} \right)^\top = \mathbf{S}_d^{(t)}(\mathbf{x}^{(t)}) \mathbf{W}_d^{(t)} \dots \mathbf{S}_{l+1}^{(t)}(\mathbf{x}^{(t)}) \mathbf{W}_{l+1}^{(t)} \mathbf{S}_l^{(t)}(\mathbf{x}^{(t)}).$$

Note that the above equation abuses notation for the last layer activation  $\mathbf{S}_d^{(t)}(\mathbf{x}^{(t)})$ . Specifically, depending on whether we are using an identity or sigmoid activation function in the last layer, we set

$$\mathbf{S}_d^{(t)}(\mathbf{x}^{(t)}) = 1 \quad \text{or} \quad \mathbf{S}_d^{(t)}(\mathbf{x}^{(t)}) = \psi\left(\mathbf{W}_d^{(t)}\left(f_{d-1}^{(t)} \circ \dots \circ f_1^{(t)}(\mathbf{x}^{(t)})\right) + \mathbf{b}_d^{(t)}\right), \quad (2)$$

where  $\psi(z) = \frac{\partial \rho_d(z)}{\partial z} = \frac{1}{1+e^{-z}} \cdot \left(1 - \frac{1}{1+e^{-z}}\right)$ .

**Part 1.** Define the following shorthand notation:

$$\left\| \left( \frac{\partial \ell(f(\mathbf{x}^{(t)}, \mathbf{w}^{(t)}), y^{(t)})}{\partial \mathbf{b}_l} \right)^\top \right\|_2 = \beta_l(\mathbf{x}^{(t)})$$

It follows from definition that

$$\begin{aligned} \beta_l(\mathbf{x}^{(t)}) &= \left\| \frac{\partial \ell(o, y^{(t)})}{\partial o} \mathbf{S}_d^{(t)}(\mathbf{x}^{(t)}) \mathbf{W}_d^{(t)} \dots \mathbf{S}_l^{(t)}(\mathbf{x}^{(t)}) \mathbf{W}_{l+1}^{(t)} \mathbf{S}_l(\mathbf{x}^{(t)}) \right\|_2 \\ &= \left| \frac{\partial \ell(o, y^{(t)})}{\partial o} \right| \left\| \left( \prod_{i=d}^{l+1} \mathbf{S}_i^{(t)}(\mathbf{x}^{(t)}) \mathbf{W}_i^{(t)} \right) \mathbf{S}_l^{(t)}(\mathbf{x}^{(t)}) \right\|_2 \end{aligned}$$

or equivalently,

$$\left\| \left( \prod_{i=d}^{l+1} \mathbf{S}_i^{(t)}(\mathbf{x}^{(t)}) \mathbf{W}_i^{(t)} \right) \mathbf{S}_l^{(t)}(\mathbf{x}^{(t)}) \right\|_2 = \beta_l(\mathbf{x}^{(t)}) \left| \frac{\partial \ell(o, y^{(t)})}{\partial o} \right|^{-1} = \frac{\beta_l(\mathbf{x}^{(t)})}{\epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})}. \quad (3)$$

**Part 2.** We are interested in upper bounding the Lipschitz constant  $\lambda_{f^{(t)}}$  of the NN close to the training data  $X$ .

First observe that  $f(\mathbf{x}, \mathbf{w}) = f_{d \leftarrow 2} \circ f_1(\mathbf{x}, \mathbf{w})$ , where we set

$$f_{d \leftarrow 2}(\mathbf{x}, \mathbf{w}) = f_d \circ f_{d-1} \circ \dots \circ f_2(\mathbf{x}, \mathbf{w})$$

Let  $\mathcal{R}_{\mathbf{x}^{(t)}}$  be the region associated with point  $\mathbf{x}^{(t)}$  and  $\mathcal{R}_{f_1^{(t)}(\mathbf{x}^{(t)})}$  the region of the NN  $f_{d \leftarrow 2}^{(t)}$  associated with point  $f_1^{(t)}(\mathbf{x}^{(t)})$ . The Lipschitz constants of  $f_{d \leftarrow 2}^{(t)}$  and  $f^{(t)}$  are related as follows:

$$\lambda_{f_{d \leftarrow 2}^{(t)}}(\mathcal{R}_{f_1^{(t)}(\mathbf{x}^{(t)})}) \cdot \sigma_n(\mathbf{W}_1^{(t)}) \leq \lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}^{(t)}}) \leq \lambda_{f_{d \leftarrow 2}^{(t)}}(\mathcal{R}_{f_1^{(t)}(\mathbf{x}^{(t)})}) \cdot \sigma_1(\mathbf{W}_1^{(t)}), \quad (4)$$

whereas

$$\lambda_{f_{d \leftarrow 2}^{(t)}}(\mathcal{R}_{f_1^{(t)}(\mathbf{x}^{(t)})}) = \left\| \left( \prod_{l=d}^2 \mathbf{S}_l^{(t)}(\mathbf{x}^{(t)}) \mathbf{W}_l^{(t)} \right) \mathbf{S}_1^{(t)}(\mathbf{x}^{(t)}) \right\|_2 = \frac{\beta_1(\mathbf{x}^{(t)})}{\epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})}. \quad (5)$$

Combining (4) with (5), we obtain

$$\frac{\beta_1(\mathbf{x}^{(t)})}{\epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \cdot \sigma_n(\mathbf{W}_1^{(t)}) \leq \lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}^{(t)}}) \leq \frac{\beta_1(\mathbf{x}^{(t)})}{\epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \cdot \sigma_1(\mathbf{W}_1^{(t)})$$

**Part 3.** Re-organizing the SGD expression and taking the norm, we have

$$\beta_l(\mathbf{x}^{(t)}) = \left\| \left( \frac{\partial \ell(f(\mathbf{x}^{(t)}, \mathbf{w}^{(t-1)}), y^{(t)})}{\partial \mathbf{b}_l^{(t)}} \right)^\top \right\|_2 = \frac{1}{\alpha_t} \|\mathbf{b}_l^{(t+1)} - \mathbf{b}_l^{(t)}\|_2.$$

implying also

$$\frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\alpha_t \cdot \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \cdot \sigma_n(\mathbf{W}_1^{(t)}) \leq \lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}^{(t)}}) \leq \frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\alpha_t \cdot \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \cdot \sigma_1(\mathbf{W}_1^{(t)}),$$

as claimed.

## B.2 Proof of Theorem 1

The proof of the theorem follows directly from Lemma 1 by summing over the training trajectory:

$$\sum_{t \in T} \frac{\lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}})}{\sigma_n(\mathbf{W}_1^{(t)})} \stackrel{\text{Lemma 1}}{\geq} \sum_{t \in T} \frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\alpha_t \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \stackrel{\text{Lemma 1}}{\geq} \sum_{t \in T} \frac{\lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}})}{\sigma_1(\mathbf{W}_1^{(t)})}.$$

## B.3 Proof of Corollary 1

For any point  $\mathbf{x}$  with label  $y$  and iteration  $t \in T$ , we say that condition  $c_t(\mathbf{x})$  holds if

$$\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\| \leq \varphi \alpha_t \epsilon_{f^{(t)}}(\mathbf{x}, y)$$

The above is the same condition as in the corollary statement but applied to an arbitrary point  $\mathbf{x}$ .

Write  $\kappa_t$  to refer to the number of training points  $\mathbf{x}_i \in X$  for which  $c_t(\mathbf{x}_i)$  holds: clearly,  $\kappa_t \in [0, N]$ , where  $N$  is the size of the training set.

We also suppose that there are  $\xi$  iterations within  $T$  for which  $\kappa_t < N$ : these are the iterations where the Lipschitz constant is larger than  $\beta\varphi$  for at least one point in the training set.

Since we sample  $\mathbf{x}^{(t)}$  i.i.d. from  $X$ , the probability that  $c_t(\mathbf{x}^{(t)})$  is satisfied for every  $t \in T$  is at most

$$\prod_{t \in T} \frac{\kappa_t}{N} \leq \left(\frac{N-1}{N}\right)^\xi = \left(1 - \frac{1}{N}\right)^\xi \leq e^{-\xi/N}.$$

By noting that above corresponds to the probability that the NN is  $(\tau, \varphi)$ -steady, we deduce that  $\xi$  cannot grow with  $|T|$  (otherwise, the probability that the NN is  $(\tau, \varphi)$ -steady would become zero as  $|T| \rightarrow \infty$ , which contradicts the corollary assumptions).

To complete the derivation, we note that if we select the iteration  $t$  at random from  $T$ , the probability that there will be some  $\mathbf{x}_i \in X$  for which  $c_t(\mathbf{x}_i)$  is not satisfied is  $\xi/|T| = O(1/|T|)$ , which converges to 0 as  $|T|$  grows.

## B.4 Proof of Corollary 2

We consider the interval  $T = \{t_1 + 1, \dots, t_2\}$  and fix

$$\mathbf{b}_1 = \operatorname{argmin}_{\mathbf{b} \in \mathbb{R}^n} \sum_{t \in T} \|\mathbf{b}_1^{(t)} - \mathbf{b}\|_2^2 = \sum_{t \in T} \frac{\mathbf{b}_1^{(t)}}{|T|}$$

to be the average bias. Working as in the proof of Theorem 1, we deduce

$$\sum_{t \in T} \frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \geq \sum_{t \in T} \frac{\alpha_t \lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}^{(t)}})}{\sigma_1(\mathbf{W}_1^{(t)})}.$$

We then proceed to upper bound the trajectory length studied in Theorem 1 in terms of the (empirical) variance of the bias:

$$\begin{aligned} \left( \sum_{t \in T} \frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \right)^2 &\leq \left( \sum_{t \in T} \frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1\|_2 + \|\mathbf{b}_1^{(t)} - \mathbf{b}_1\|_2}{\epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \right)^2 \\ &\leq \left( \sum_{t \in T} \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})^{-2} \right) \left( \sum_{t \in T} \left( \|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1\|_2 + \|\mathbf{b}_1^{(t)} - \mathbf{b}_1\|_2 \right)^2 \right) \\ &\hspace{15em} \text{(From Cauchy's inequality)} \\ &\leq \sum_{t \in T} \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})^{-2} \sum_{t \in T} 2 \left( \|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1\|_2^2 + \|\mathbf{b}_1^{(t)} - \mathbf{b}_1\|_2^2 \right) \\ &\hspace{15em} \text{(since } (a+b)^2 \leq 2(a^2 + b^2)\text{)} \\ &\leq 4 \sum_{t \in T} \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})^{-2} \sum_{t=t_1}^{t_2} \|\mathbf{b}_1^{(t)} - \mathbf{b}_1\|_2^2 \end{aligned}$$

or, equivalently,

$$\sum_{t=t_1}^{t_2} \frac{\|\mathbf{b}_1^{(t)} - \mathbf{b}_1\|_2^2}{|T|} \geq \frac{1}{4|T|} \left( \sum_{t \in T} \frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \right)^2 \frac{1}{\sum_{t \in T} \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})^{-2}} \quad (6)$$

Thus, we have

$$\sum_{t=t_1}^{t_2} \frac{\|\mathbf{b}_1^{(t)} - \mathbf{b}_1\|_2^2}{|T|} \geq 0.25 \left( \underset{t=t_1}{\text{avg}} \frac{\alpha_t \lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}^{(t)}})}{\sigma_1(\mathbf{W}_1^{(t)})} \right)^2 \frac{|T|}{\sum_{t \in T} \frac{1}{\epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})^2}}. \quad (7)$$

The proof concludes by noticing that the right-most term corresponds to a harmonic mean of  $\epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})^2$  over  $t \in T$ .

### B.5 Proof of Corollary 3

Suppose that we train our NN for  $\tau$  iterations and set  $T = \{0, \dots, \tau-1\}$ . The distance to initialization is bounded by

$$\|\mathbf{b}_1^{(\tau)} - \mathbf{b}_1^{(0)}\|_2 \leq \sum_{t \in T} \|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2 = \sum_{t \in T} \frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)}).$$

We obtain the final expression by arguing as in the proof of Theorem 1 to write:

$$\sum_{t \in T} \frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)}) \leq \sum_{t \in T} \frac{\alpha_t \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)}) \lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}^{(t)}})}{\sigma_n(\mathbf{W}_1^{(t)})}.$$

### B.6 Proof of Theorem 2

We will start by proving the following Lemma:

**Lemma 2.** *Let  $\mathcal{R}_{\mathbf{x}}$  be a linear region of  $f^{(t)}$  and suppose that there exists a vector  $\mathbf{a} \in \mathbb{R}^{|T|}$  such that*

$$\prod_{l=d-1}^1 [\mathcal{S}_l^{(t)}(\mathbf{x})](i_l, i_l) = \sum_{k \in T} \mathbf{a}(k) \cdot \prod_{l=d-1}^1 [\mathcal{S}_l^{(k)}(\mathbf{x}^{(k)})](i_l, i_l) \quad (8)$$

for all indices  $\{i_l\}_{l=1, \dots, d-1}$ , with  $i_l \in \{1, \dots, n_l\}$ . Then,  $f^{(t)}$  is Lipschitz continuous within  $\mathcal{R}_{\mathbf{x}}$  and its Lipschitz constant is at most

$$\lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}}) \leq (1 + \gamma) \sum_{k \in T} |\mathbf{a}(k)| \frac{|\mathcal{S}_d^{(t)}(\mathbf{x})|}{|\mathcal{S}_d^{(k)}(\mathbf{x}^{(k)})|} \lambda_{f^{(k)}}(\mathcal{R}_{\mathbf{x}^{(k)}}),$$

with

$$\mathcal{S}_d^{(t)}(\mathbf{z}) = \frac{\partial \rho_d(\mathbf{x})}{\partial \mathbf{x}} \Big|_{\mathbf{x} = \mathbf{W}_d^{(t)}(f_{d-1}^{(t)} \circ \dots \circ f_1^{(t)}(\mathbf{z}) + \mathbf{b}_d^{(t)}}. \quad (9)$$

*Proof.* The gradient of a network  $f^{(t)}$  at point  $\mathbf{x}$  is simply

$$\nabla f^{(t)}(\mathbf{x}) = \prod_{l=d}^1 \mathcal{S}_l^{(t)}(\mathbf{x}) \mathbf{W}_l^{(t)} = \mathcal{S}_d^{(t)}(\mathbf{x}) \mathbf{W}_d^{(t)} \overbrace{\prod_{l=d-1}^1 \mathcal{S}_l^{(t)}(\mathbf{x}) \mathbf{W}_l^{(t)}}^{q^{(t)}(\mathbf{x})}.$$

Term  $q^{(t)}(\mathbf{x})$  can be expanded as follows:

$$q^{(t)}(\mathbf{x}) = \sum_{i_{d-1}=1}^{n_d} \mathbf{W}_d^{(t)}(1, i_{d-1}) \left[ \prod_{l=d-1}^1 \mathcal{S}_l^{(t)}(\mathbf{x}) \mathbf{W}_l^{(t)} \right] (i_{d-1}, :)$$



$$\begin{aligned}
&= \dots \\
&= \sum_{i_{d-1}=1}^{n_{d-1}} \dots \sum_{i_1=1}^{n_1} \mathbf{W}_d^{(t)}(1, i_{d-1}) \dots [\mathbf{S}_1^{(t)}(\mathbf{x})](i_1, i_1) \mathbf{W}_1^{(t)}(i_1, :).
\end{aligned}$$

Thus, under the main Lemma condition it is also true that

$$\begin{aligned}
q^{(t)}(\mathbf{x}) &= \sum_{k \in T} \mathbf{a}(k) \cdot \left( \sum_{i_{d-1}=1}^{n_{d-1}} \dots \sum_{i_1=1}^{n_1} \mathbf{W}_d^{(t)}(1, i_{d-1}) \dots [\mathbf{S}_1^{(k)}(\mathbf{x}^{(k)})](i_1, i_1) \mathbf{W}_1^{(t)}(i_1, :) \right) \\
&= \sum_{k \in T} \mathbf{a}(k) \cdot \left( \sum_{i_{d-1}=1}^{n_{d-1}} \dots \sum_{i_1=1}^{n_1} \mathbf{W}_d^{(t)}(1, i_{d-1}) \dots [\mathbf{S}_1^{(t)}(\mathbf{x}^{(k)})](i_1, i_1) \mathbf{W}_1^{(t)}(i_1, :) \right) \\
&= \sum_{k \in T} \mathbf{a}(k) \cdot q^{(t)}(\mathbf{x}^{(k)}),
\end{aligned}$$

Note that in the second step above we have used Assumption 1 to argue that the training point activation patterns do not change within  $T$ .

The above analysis implies that

$$\begin{aligned}
\lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}}) &= |\mathbf{S}_d^{(t)}(\mathbf{x})| \|q^{(t)}(\mathbf{x})\|_2 = |\mathbf{S}_d^{(t)}(\mathbf{x})| \left\| \sum_{k \in T} \mathbf{a}(k) \cdot q^{(t)}(\mathbf{x}^{(k)}) \right\|_2 \\
&\leq |\mathbf{S}_d^{(t)}(\mathbf{x})| \sum_{k \in T} \|\mathbf{a}(k) \cdot q^{(t)}(\mathbf{x}^{(k)})\|_2 \\
&= \sum_{k \in T} |\mathbf{a}(k)| \cdot \frac{|\mathbf{S}_d^{(t)}(\mathbf{x})|}{|\mathbf{S}_d^{(t)}(\mathbf{x}^{(k)})|} \cdot \lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}^{(k)}}) \\
&\leq (1 + \gamma) \sum_{k \in T} |\mathbf{a}(k)| \cdot \frac{|\mathbf{S}_d^{(t)}(\mathbf{x})|}{|\mathbf{S}_d^{(t)}(\mathbf{x}^{(k)})|} \cdot \lambda_{f^{(k)}}(\mathcal{R}_{\mathbf{x}^{(k)}})
\end{aligned}$$

with the 3rd step being true due to the triangle inequality and the 5th follows from Assumption 1.  $\square$

The proof continues by realizing that, for every index set  $i_{d-1}, \dots, i_1$  there exists an entry  $i$  such that

$$\left[ \bigotimes_{l=d-1}^1 \mathbf{S}_l^{(t)}(\mathbf{x}) \right](i, i) = \prod_{l=d-1}^1 [\mathbf{S}_l^{(t)}(\mathbf{x})](i_l, i_l).$$

Therefore, condition (8) is equivalent to asserting that

$$\begin{aligned}
\mathbf{s}_t(\mathbf{x}) &= \bigotimes_{l=d-1}^1 \text{diag}(\mathbf{S}_l^{(t)}(\mathbf{x})) = \sum_{k \in T} \mathbf{a}(k) \cdot \bigotimes_{l=d-1}^1 \text{diag}(\mathbf{S}_l^{(k)}(\mathbf{x}^{(k)})) \\
&= \sum_{k \in T} \mathbf{a}(k) \cdot \mathbf{s}_k(\mathbf{x}^{(k)}) = \mathbf{S}_T \mathbf{a}.
\end{aligned}$$

Let us focus on  $|\mathbf{S}_d^{(t)}(\mathbf{x})|/|\mathbf{S}_d^{(t)}(\mathbf{x}^{(k)})|$ . When there is no activation in the last layer, the term is trivially  $\xi = 1$ . We next derive an upper bound to also account for the sigmoid activation: To do this, set  $z = \mathbf{W}_d^{(t)}(f_{d-1}^{(t)} \circ \dots \circ f_1^{(t)}(\mathbf{x})) + \mathbf{b}_d^{(t)}$  such that

$$\mathbf{S}_d^{(t)}(\mathbf{x}) = \psi(z) \quad \text{with} \quad \psi(z) = \frac{\partial \rho_d(z)}{\partial z} = \frac{1}{1 + e^{-z}} \cdot \left( 1 - \frac{1}{1 + e^{-z}} \right)$$

Function  $\psi$  takes its maximum value for  $z = 0$ , with  $\psi(z) \leq \psi(0) = 0.25$ . We notice that  $\psi$  is symmetric around 0 and monotonically decreasing on either side. Its minimum is thus given when  $|z|$  is as large as possible. However, since our classifier's output is bounded in  $f^{(t)}(\mathbf{x}) \in [\mu_T, 1 - \mu_T]$  for all points seen within  $T$ , we have  $|z| \leq \log(1/\mu_T - 1)$  and thus

$$|\mathbf{S}_d^{(t)}(\mathbf{x}^{(k)})| \geq \psi(\log(1/\mu_T - 1)) = \mu_T(1 - \mu_T).$$

All in all, we get  $|\mathbf{S}_d^{(t)}(\mathbf{x})|/|\mathbf{S}_d^{(t)}(\mathbf{x}^{(k)})| \leq 0.25/(\mu_T(1 - \mu_T)) = \xi$ .

We then rely on Lemma 1 to upper bound each local Lipschitz constant in terms of the bias update:

$$\lambda_{f^{(k)}}(\mathcal{R}_{\mathbf{x}^{(k)}}) \leq \frac{\|\mathbf{b}_1^{(k+1)} - \mathbf{b}_1^{(k)}\|_2}{\alpha_k \epsilon_{f^{(k)}}(\mathbf{x}^{(k)}, \mathbf{y}^{(k)})} \sigma_1(\mathbf{W}_1^{(k)})_2 \leq \beta \frac{\|\mathbf{b}_1^{(k+1)} - \mathbf{b}_1^{(k)}\|_2}{\alpha_k \epsilon_{f^{(k)}}(\mathbf{x}^{(k)}, \mathbf{y}^{(k)})}, \quad (10)$$

matching the claim of the theorem.

### B.7 Proof of Theorem 3

We repeat the theorem statement here for easy reference:

**Theorem 3.** *Let  $f^{(t)}$  be a depth  $d$  NN with ReLU activations being trained with SGD, a BCE loss and  $1/2$ -Dropout.*

*Suppose that  $f^{(t)}$  is  $(\tau, \varphi)$ -steady and that for every  $t \geq \tau$  the following hold: (a) Assumption 1, (b)  $\mathbf{s}_t(\mathbf{x}) \leq \sum_{i=1}^N \mathbf{s}_t(\mathbf{x}_i)$  for every  $\mathbf{x} \in \mathcal{X}$ , (c)  $\sigma_1(\mathbf{W}_1^{(t)}) \leq \beta$ , and (d)  $f^{(t)}(\mathbf{x}^{(t)}) \in [\mu, 1 - \mu]$ .*

*Define*

$$r_t(X) = \frac{\min_{i=1}^N |1 - 2f^{(t)}(\mathbf{x}_i)|}{c \varphi \log\left(\sum_{l=1}^{d-1} n_l\right)} \quad \text{and} \quad c = \frac{(1 + \gamma) \beta (1 + o(1))}{\mu (1 - \mu) p_{\min}},$$

*where  $p_{\min} = \min_{l < d, i \leq n_l, t \geq \tau} [\text{avg}_{\mathbf{x} \in X} \text{diag}(\mathbf{S}_l^{(t)}(\mathbf{x}))]_i > 0$  is the minimum frequency that any neuron is active before Dropout is applied.*

*For any  $\delta > 0$ , with probability at least  $1 - \delta$  over the Dropout and the training set sampling, the generalization error is at most*

$$|\text{er}_t^{\text{emp}} - \text{er}_t^{\text{exp}}| \leq \sqrt{\frac{4 \log(2) \mathcal{N}(\mathcal{X}; \ell_2, r(X)) + 2 \log(1/\delta)}{N}},$$

*where  $\mathcal{N}(\mathcal{X}; \ell_2, r)$  is the minimal number of  $\ell_2$ -balls of radius  $r$  needed to cover  $\mathcal{X}$ .*

The proof consists of two parts. First, Lemma 3 provides a bound on the *global* Lipschitz constant of a NN trained with Dropout as a function of the bias updates observed during a sufficiently long training. Then, Lemma 4 uses techniques from the robustness framework [2, 3] to derive a generalization bound.

#### B.7.1 The global Lipschitz constant

We prove the following:

**Lemma 3.** *In the setting of Theorem 2, suppose that the network is trained using  $1/2$ -Dropout and denote by  $\mathbf{p}_l = \text{avg}_{\mathbf{x} \in X} \text{diag}(\mathbf{S}_l^{(t)}(\mathbf{x}))$  the probability that the neurons in layer  $l$  are active (before Dropout is applied). The global Lipschitz constant of  $f^{(t)}$  is with high probability*

$$\lambda_{f^{(t)}} \leq c \log\left(\sum_{l=1}^{d-1} n_l\right) \|\varphi_T\|_{\infty} := \lambda_{f^{(t)}}^{\text{steady}}$$

*for  $c = (1 + \gamma) \beta \xi (1 + o(1)) / p_{\min}$ , whenever  $|T| = \tilde{\Omega}\left(\frac{p_{\text{avg}}}{p_{\min}^2} \sum_{l=1}^{d-1} n_l\right)$ , with  $p_{\min}$  and  $p_{\text{avg}}$  being the minimum and average entry over all  $\mathbf{p}_l$ , respectively.*

The inequality provided above is unexpectedly tight: combining  $\lambda_{f^{(t)}} \geq \lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}^{(t)}})$  with Lemma 1 we can deduce that

$$\lambda_{f^{(t)}} \leq \lambda_{f^{(t)}}^{\text{steady}} \leq \lambda_{f^{(t)}} O(\log(dn)),$$

where we have assumed that  $c/\sigma_n(\mathbf{W}_1) = O(1)$  and  $n_l = n$  for all  $l < d$ .

*Proof.* Let  $\mathbf{x}$  be a point within a region where  $f^{(t)}$  assumes its maximum gradient norm.

The activation  $\tilde{\mathbf{S}}_T(:, t) = \tilde{\mathbf{s}}_t(\mathbf{x}^{(t)})$  at the  $t$ -th SGD iteration is obtained by a two step procedure:

1. Sample a point  $\mathbf{x}^{(t)}$  from  $X$  with replacement. Let  $\mathbf{s}_t(\mathbf{x}^{(t)}) = \bigotimes_{l=d-1}^1 \mathbf{s}_{t,l}(\mathbf{x}^{(t)})$  be its activation pattern (before dropout), where  $\mathbf{s}_{t,l}(\mathbf{x}^{(t)}) := \text{diag}\left(\mathbf{S}_l^{(t)}(\mathbf{x}^{(t)})\right)$ .
2. Construct  $\tilde{\mathbf{s}}_t(\mathbf{x}^{(t)})$  by setting each neuron activation to zero with probability 0.5. Specifically,  $\tilde{\mathbf{s}}_t(\mathbf{x}^{(t)}) = \bigotimes_{l=d-1}^1 (\mathbf{z}_l \circ \mathbf{s}_{t,l}(\mathbf{x}^{(t)}))$ , where  $\mathbf{z}_l \in \{0, 1\}^{n_l}$  is a random binary vector.

Let  $\mathbf{S}$  be a binary matrix containing neuron activations as columns. We introduce the following definitions:

- We call  $\mathbf{S}$  a *covering set* if  $\mathbf{S}\mathbf{1} \geq \mathbf{1}$  with the inequality taken element-wise.
- We call  $\mathbf{S}$  a *basis* of  $\mathbf{s}_t(\mathbf{x})$  if  $\mathbf{S}\mathbf{1} = \mathbf{s}_t(\mathbf{x})$ .

Our proof hinges on two observations:

*Observation 1.* Every basis yields a bound on the Lipschitz constant of  $f^{(t)}$  (this can be seen from the proof of Theorem 2). Specifically, for any  $k$  training points  $\mathbf{x}_1, \dots, \mathbf{x}_k$  whose activations  $\mathbf{S} = [\mathbf{s}_t(\mathbf{x}_1), \dots, \mathbf{s}_t(\mathbf{x}_k)]$  is a basis of  $\mathbf{s}_t(\mathbf{x})$ , we have

$$\lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}}) \leq \xi \sum_{i=1}^k \lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}_i}) \leq k \xi \max_{i=1}^k \lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}_i}).$$

where  $\xi \geq \frac{|\mathbf{S}_d^{(t)}(\mathbf{x})|}{|\mathbf{S}_d^{(t)}(\mathbf{x}_i)|}$  accounts for the sigmoid.

Thus, if we don't use dropout and within the columns of  $\mathbf{S}_T$  there exist  $k$  that form a basis of  $\mathbf{s}_t(\mathbf{x})$ , then this also implies that the global Lipschitz constant will be bounded by

$$\lambda_{f^{(t)}}(\mathcal{R}_{\mathbf{x}}) \leq k\beta(1 + \gamma)\xi \max_{t \in T} \frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\alpha_{t \in f^{(t)}}(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})},$$

where, in an identical fashion to Theorem 2, the  $1 + \gamma$  factor is added due to Assumption 1 in order to account for  $f^{(t)}$  not having completely converged, and we have also used Lemma 1 and the uniform bound  $\|\mathbf{W}_1^{(t)}\|_2 \leq \beta$ .

*Observation 2.* Let us consider the effect of Dropout. Suppose that  $\mathbf{S}_T$  does not contain a basis of  $\mathbf{s}_t(\mathbf{x}) = \bigotimes_{l=d-1}^1 \mathbf{s}_{t,l}(\mathbf{x})$ , but there exist a set of columns  $\mathbf{S}$  that is a covering set (as we will see, this is a much easier condition to satisfy). Denote by  $\tilde{\mathbf{S}}$  the same matrix after the Dropout sampling. Then, with some strictly positive probability,  $\tilde{\mathbf{S}}$  can become a basis.

**Claim 1.** For any  $k$  training points  $\mathbf{x}_1, \dots, \mathbf{x}_k$  whose activations  $\mathbf{S} = [\mathbf{s}_t(\mathbf{x}_1), \dots, \mathbf{s}_t(\mathbf{x}_k)]$  form a covering set, there must exist  $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_k]$  with  $\mathbf{q}_i = \bigotimes_{l=d-1}^1 \mathbf{q}_{i,l}$  and  $\mathbf{q}_{i,l} \leq \mathbf{s}_{t,l}(\mathbf{x}_i)$  (i.e., that Dropout can sample) such that  $\mathbf{Q}$  is a basis of  $\mathbf{s}_t(\mathbf{x})$ .

*Proof.* To deduce this fact, we notice that since

$$\sum_{i=1}^k \mathbf{q}_i = \sum_{i=1}^k \bigotimes_{l=d-1}^1 \mathbf{q}_{i,l} = \bigotimes_{l=d-1}^1 \left( \sum_{i=1}^k \mathbf{q}_{i,l} \right) \quad \text{and} \quad \mathbf{s}_t(\mathbf{x}) = \bigotimes_{l=d-1}^1 \mathbf{s}_{t,l}(\mathbf{x}),$$

to ensure that  $\mathbf{Q}$  is a basis we need to show that, for every  $l$ , there exists  $[\mathbf{q}_{1,l} \cdots \mathbf{q}_{k,l}]$  with  $\mathbf{q}_{i,l} \leq \mathbf{s}_{t,l}(\mathbf{x}_i)$  such that  $\sum_{i=1}^k \mathbf{q}_{i,l} = \mathbf{s}_{t,l}(\mathbf{x})$ . The latter can always be satisfied when  $\sum_{i=1}^k \mathbf{s}_{t,l}(\mathbf{x}_i) \geq \mathbf{1}$ . When  $\mathbf{S}$  is a covering set we have

$$\mathbf{S}\mathbf{1} = \sum_{i=1}^k \bigotimes_{l=d-1}^1 \mathbf{s}_{t,l}(\mathbf{x}_i) = \bigotimes_{l=d-1}^1 \sum_{i=1}^k \mathbf{s}_{t,l}(\mathbf{x}_i) \geq \mathbf{1},$$

which also implies  $\sum_{i=1}^k \mathbf{s}_{t,l}(\mathbf{x}_i) \geq \mathbf{1}$  as needed.  $\square$

To obtain an upper bound for the Lipschitz constant of  $f^{(t)}$ , our strategy will entail lower bounding the probability that such a basis of  $\mathbf{s}_t(\mathbf{x})$  will be seen within  $T$ .

Consider any  $k$  training points  $\mathbf{x}_1, \dots, \mathbf{x}_k$  sampled with replacement from  $X$  and let  $\mathbf{S} = [\mathbf{s}_t(\mathbf{x}_1), \dots, \mathbf{s}_t(\mathbf{x}_k)]$  be the corresponding (random) matrix of neural activations. Further, denote by  $p_{\text{cover}}(\mathbf{S})$  the probability that  $\mathbf{S}$  is a covering set.

The probability  $p_{\text{basis}}(\mathbf{S}_T)$  that  $\tilde{\mathbf{S}}_T$  contains a basis of  $\mathbf{s}_t(\mathbf{x})$  is given by

$$\begin{aligned} p_{\text{basis}}(\mathbf{S}_T) &= 1 - \mathbb{P}\left(\tilde{\mathbf{S}}_T \text{ does not contain a basis}\right) \\ &\geq 1 - \prod_{p=1}^{\lfloor \frac{|T|}{k} \rfloor} \mathbb{P}\left(\tilde{\mathbf{S}}_T(:, (p-1)k+1 : pk) \text{ is not a basis}\right) \end{aligned}$$

For every  $\tilde{\mathbf{S}}_T(:, (p-1)k+1 : pk)$  we have:

$$\begin{aligned} \mathbb{P}\left(\tilde{\mathbf{S}}_T(:, (p-1)k+1 : pk) \text{ is a basis}\right) &= \mathbb{P}\left(\tilde{\mathbf{S}} \text{ is a basis}\right) \\ &= \mathbb{P}\left(\tilde{\mathbf{S}} \text{ is a basis} \mid \mathbf{S} \text{ is a covering set}\right) \mathbb{P}(\mathbf{S} \text{ is a covering set}) \\ &= \mathbb{P}\left(\tilde{\mathbf{S}} \text{ is a basis} \mid \mathbf{S} \text{ is a covering set}\right) p_{\text{cover}}(\mathbf{S}). \end{aligned}$$

By Observation 2, if  $\mathbf{S}$  is a covering set then there must exist  $\mathbf{q}_i = \bigotimes_{l=d-1}^1 \mathbf{q}_{t,l}(\mathbf{x}_i) \leq \mathbf{s}_t(\mathbf{x}_i) = \bigotimes_{l=d-1}^1 \mathbf{s}_{t,l}(\mathbf{x}_i)$ , such that  $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_k]$  is a basis of  $\mathbf{s}_t(\mathbf{x})$ .

We proceed to compute the probability that the activation pattern sampled by Dropout  $\tilde{\mathbf{s}}_t(\mathbf{x}_i) = \bigotimes_{l=d}^2 (\mathbf{z}_{i,l} \circ \mathbf{s}_{t,l}(\mathbf{x}_i))$ , where  $\mathbf{z}_{i,t}$  are random binary vectors, is a basis of  $\mathbf{s}_t(\mathbf{x}^{(t)})$  due to  $\tilde{\mathbf{S}} = \mathbf{Q}$ :

$$\begin{aligned} \mathbb{P}\left(\tilde{\mathbf{S}} \text{ is a basis} \mid \mathbf{S} \text{ is a covering set}\right) &= \mathbb{P}\left(\tilde{\mathbf{S}} = \mathbf{Q}\right) \\ &= \mathbb{P}(\tilde{\mathbf{s}}_t(\mathbf{x}_i) = \mathbf{q}_i \text{ for } i = 1, \dots, k) \\ &= \prod_{i=1}^k \mathbb{P}(\tilde{\mathbf{s}}_t(\mathbf{x}_i) = \mathbf{q}_i) \\ &= \prod_{i=1}^k \mathbb{P}\left(\bigotimes_{l=d-1}^1 (\mathbf{z}_{i,l} \circ \mathbf{s}_{t,l}(\mathbf{x}_i)) = \bigotimes_{l=d-1}^1 \mathbf{q}_{i,l}\right) \\ &= \prod_{i=1}^k \prod_{l=d-1}^1 \mathbb{P}(\mathbf{z}_{i,l} \circ \mathbf{s}_{t,l}(\mathbf{x}_i) = \mathbf{q}_{i,l}) \\ &= \prod_{i=1}^k \prod_{l=d-1}^1 \frac{1}{2^{\|\mathbf{s}_{t,l}(\mathbf{x}_i)\|_1}} = 2^{-\sum_{i=1}^k \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x}_i)\|_1}, \end{aligned}$$

where the second to last step is a consequence of Dropout with probability 0.5 sampling for each layer uniformly at random from the set of all possible neuron activation patterns that can be obtained by disabling some neurons of  $\mathbf{s}_{t,l}(\mathbf{x}_i)$ .

Term  $\sum_{i=1}^k \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x}_i)\|_1$  can be seen as the sum of  $k$  independent samples, each having mean  $m = \text{avg}_{\mathbf{x} \in X} \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x})\|_1$  and maximum value  $c = \max_{\mathbf{x} \in X} \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x})\|_1$ . Hoeffding's inequality yields

$$\mathbb{P}\left(\sum_{i=1}^k \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x}_i)\|_1 > \mathbb{E}\left[\sum_{i=1}^k \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x}_i)\|_1\right] + k\delta\right) < \exp\left(-\frac{2k^2\delta^2}{kc^2}\right),$$

implying also that  $\mathbb{P}\left(2^{-\sum_{i=1}^k \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x}_i)\|_1} < 2^{-(k\mu + \sqrt{k/2c})}\right) < 1/e$ . Thus,

$$\begin{aligned} &\mathbb{P}\left(\tilde{\mathbf{S}} \text{ is a basis} \mid \mathbf{S} \text{ is a covering set}\right) \\ &= \mathbb{P}\left(\tilde{\mathbf{S}} \text{ is a basis} \mid \mathbf{S} \text{ is a covering set}, 2^{-\sum_{i=1}^k \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x}_i)\|_1} \geq 2^{-h}\right) \mathbb{P}\left(2^{-\sum_{i=1}^k \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x}_i)\|_1} < 2^{-h}\right) \end{aligned}$$

$$\begin{aligned}
& + \mathbf{P}\left(\tilde{\mathcal{S}} \text{ is a basis} \mid \mathcal{S} \text{ is a covering set}, 2^{-\sum_{i=1}^k \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x}_i)\|_1} < 2^{-h}\right) \mathbf{P}\left(2^{-\sum_{i=1}^k \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x}_i)\|_1} < 2^{-h}\right) \\
& \geq \mathbf{P}\left(\tilde{\mathcal{S}} \text{ is a basis} \mid \mathcal{S} \text{ is a covering set}, 2^{-\sum_{i=1}^k \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x}_i)\|_1} \geq 2^{-h}\right) \mathbf{P}\left(2^{-\sum_{i=1}^k \sum_{l=d-1}^1 \|\mathbf{s}_{t,l}(\mathbf{x}_i)\|_1} < 2^{-h}\right) \\
& \geq 2^{-(k\mu+c\sqrt{k/2})} (1-1/e) > 2^{-(k\mu+c\sqrt{k/2+1})},
\end{aligned}$$

where the first step employs the law of total probability. We therefore deduce that

$$\begin{aligned}
p_{\text{basis}}(\mathcal{S}_T) & \geq 1 - \left(\frac{p_{\text{cover}}(\mathcal{S})}{2^{(k\mu+c\sqrt{k/2+1})}}\right)^{\lfloor \frac{|T|}{k} \rfloor} \\
& = 1 - 2^{\left(\lfloor \frac{|T|}{k} \rfloor \log_2\left(\frac{p_{\text{cover}}(\mathcal{S})}{2^{(k\mu+c\sqrt{k/2+1})}}\right)\right)} \\
& = 1 - 2^{-\left(\frac{\lfloor \frac{|T|}{k} \rfloor}{(k\mu+c\sqrt{k/2+1})} \log_2(1/p_{\text{cover}}(\mathcal{S}))\right)},
\end{aligned}$$

which is satisfied with high probability when

$$|T| = \Omega\left(\frac{k(k\mu + c\sqrt{k/2} + 1)}{\log(1/p_{\text{cover}}(\mathcal{S}))}\right) = \Omega\left(\frac{k^2\mu + nk^{3/2}}{-\log p_{\text{cover}}(\mathcal{S})}\right).$$

The final step of the proof entails bounding  $\mu$  and  $p_{\text{cover}}(\mathcal{S})$ . We will think of neuron  $i$  at layer  $l$  as a (dependent) Bernoulli random variable with activation probability  $\mathbf{p}_l(i)$ . The probability that neuron  $i$  in layer  $l$  is not activated within  $k$  independent trials is  $(1 - \mathbf{p}_l(i))^k$ . Taking a union bound over all neurons in all layers, results in:

$$\begin{aligned}
p_{\text{cover}}(\mathcal{S}) & \geq 1 - \sum_{l=1}^{d-1} \sum_{i=1}^{n_l} (1 - \mathbf{p}_l(i))^k = 1 - \sum_{l=1}^{d-1} \sum_{i=1}^{n_l} \left(1 - \frac{k \mathbf{p}_l(i)}{k}\right)^k \\
& \geq 1 - \sum_{l,i} \exp(-k \mathbf{p}_l(i)) \\
& \geq 1 - \exp\left(-k p_{\min} + \log\left(\sum_{l=1}^{d-1} n_l\right)\right)
\end{aligned}$$

with  $p_{\min} = \min_{l,i} \mathbf{p}_l(i)$ . On the other hand, the average norm is given by

$$\begin{aligned}
m & = \text{avg}_{\mathbf{x} \in X} \sum_{l=1}^{d-1} \|\mathbf{s}_{t,l}(\mathbf{x})\|_1 = \sum_{\mathbf{x} \in X} \frac{\sum_{l=1}^{d-1} \sum_{i=1}^{n_l} [\mathbf{s}_{t,l}(\mathbf{x})](i)}{N} \\
& = \sum_{l=1}^{d-1} \sum_{i=1}^{n_l} \frac{\sum_{\mathbf{x} \in X} [\mathbf{s}_{t,l}(\mathbf{x})](i)}{N} \\
& = \sum_{l=1}^{d-1} \sum_{i=1}^{n_l} \mathbf{p}_l(i) = \left(\sum_{l=1}^{d-1} n_l\right) p_{\text{avg}}.
\end{aligned}$$

The number of iterations we thus need to obtain a high probability bound is thus

$$|T| = \Omega\left(\frac{k^2 \left(\sum_{l=1}^{d-1} n_l\right) p_{\text{avg}} + \left(\sum_{l=1}^{d-1} n_l\right) k^{3/2}}{-\log\left(1 - \exp\left(-k p_{\min} + \log\left(\sum_{l=1}^{d-1} n_l\right)\right)\right)}\right).$$

If we select  $k = (1 + o(1)) \log\left(\sum_{l=1}^{d-1} n_l\right) / p_{\min}$ , we obtain

$$|T| = \Omega\left(\left(\frac{(1 + o(1)) \log\left(\sum_{l=1}^{d-1} n_l\right)}{p_{\min}}\right)^2 \left(\sum_{l=1}^{d-1} n_l\right) p_{\text{avg}} + \left(\sum_{l=1}^{d-1} n_l\right) \left(\frac{(1 + o(1)) \log\left(\sum_{l=1}^{d-1} n_l\right)}{p_{\min}}\right)^{3/2}\right)$$

$$= \tilde{\Omega} \left( \left( \frac{1}{p_{\min}} \right)^2 \left( \sum_{l=1}^{d-1} n_l \right) p_{\text{avg}} + \left( \sum_{l=1}^{d-1} n_l \right) \left( \frac{1}{p_{\min}} \right)^{3/2} \right) = \tilde{\Omega} \left( \left( \sum_{l=1}^{d-1} n_l \right) \frac{p_{\text{avg}}}{p_{\min}^2} \right),$$

where the asymptotic notation hides logarithmic factors.

The final Lipschitz constant is obtained by plugging in the bound of Observation 1 the value  $k = (1 + o(1)) \log \left( \sum_{l=1}^{d-1} n_l \right) / p_{\min}$ .

□

## B.7.2 Generalization

We prove the following:

**Lemma 4.** *In the setting of Lemma 3, suppose that the NN  $f^{(t)}$  has been trained using a BCE loss and a sigmoid activation in the last layer, let  $g^{(t)}(\mathbf{x}) = \mathbf{1}[f^{(t)}(\mathbf{x}) > 0.5] \in \{0, 1\}$  the classifier's output, and define*

$$r_t(X) := \frac{\min_{i=1}^N |f^{(t)}(\mathbf{x}_i) - 0.5|}{2\lambda_{f^{(t)}}^{\text{bound}}},$$

where  $\lambda_{f^{(t)}} \leq \lambda_{f^{(t)}}^{\text{bound}}$  with probability at least  $1 - o(1)$ . For any  $\delta > 0$ , with probability at least  $1 - \delta - o(1)$ , we have

$$\left| E_{(\mathbf{x}, y)} \left[ \text{er} \left( g^{(t)}(\mathbf{x}), y \right) \right] - \text{avg}_{i=1}^N \text{er} \left( g^{(t)}(\mathbf{x}_i), y_i \right) \right| \leq \sqrt{\frac{4 \log(2) \mathcal{N}(\mathcal{X}; \ell_2, r_t(X)) + 2 \log(1/\delta)}{N}},$$

where  $\text{er}(\hat{y}, y) = \mathbf{1}[\hat{y} \neq y]$  is the classification error and  $\mathcal{N}(\mathcal{X}; \ell_2, r)$  is the minimal number of  $\ell_2$ -balls of radius  $r$  needed to cover the input domain  $\mathcal{X}$ .

*Proof.* For convenience, we drop the iteration index.

Following Xu and Mannor [2], we define the input margin  $\gamma_i$  of classifier  $g$  at  $\mathbf{x}_i$  to be

$$\gamma_i := \sup \{ a : \forall \mathbf{x}, \|\mathbf{x} - \mathbf{x}_i\|_2 \leq a, g(\mathbf{x}) = g(\mathbf{x}_i) \},$$

which is the distance (in input space) to the classification boundary. For completeness, we also repeat the definition of a robust classifier:

**Definition 1** (Adapted from Definition 2 [2]). *Classifier  $g$  is  $(K, \epsilon)$ -robust if  $\mathcal{X} \times \mathcal{Y}$  can be partitioned into  $K$  disjoint sets, denoted as  $\mathcal{C}_{k=1}^K$ , such that  $\forall i = 1 \dots, N$ ,*

$$(\mathbf{x}_i, y_i), (\mathbf{x}, y) \in \mathcal{C}_k \implies |\text{er}(g(\mathbf{x}_i), y_i) - \text{er}(g(\mathbf{x}), y)| \leq \epsilon.$$

Denote by  $\mathbf{x}_i^*$  a point with  $\|\mathbf{x}_i^* - \mathbf{x}_i\|_2 = \gamma_i$  with  $g(\mathbf{x}_i^*) = g(\mathbf{x}_i)$  and notice that  $f(\mathbf{x}_i^*) = 0.5$  (due to the definition  $g(\mathbf{x}) = \mathbf{1}[f(\mathbf{x}) > 0.5]$ ). We use the argument of Sokolić et al. [3] and bound the input margin as follows:

$$\gamma_i \geq \frac{\|f(\mathbf{x}_i) - f(\mathbf{x}_i^*)\|_2}{\lambda_f} = \frac{\|f(\mathbf{x}_i) - 0.5\|_2}{\lambda_f} \geq \frac{\|f(\mathbf{x}_i) - 0.5\|_2}{\lambda_f^{\text{bound}}}, \quad (11)$$

with probability at least  $1 - o(1)$ . From Example 1 in [2] we then deduce that  $g$  is  $(2\mathcal{N}(2\mathcal{X}, \ell_2, r_t(X)), 0)$ -robust for

$$r_t(X) = \frac{\|f(\mathbf{x}_i) - 0.5\|_2}{2\lambda_f^{\text{bound}}} \leq \min_{i=1}^N \frac{\gamma_i}{2}.$$

Theorem 3 [2] implies that if  $g$  is  $(K, 0)$ -robust then, for any  $\delta > 0$ , the following holds:

$$\left| E_{(\mathbf{x}, y)} [\text{er}(g(\mathbf{x}), y)] - \text{avg}_{i=1}^N \text{er}(g(\mathbf{x}_i), y_i) \right| \leq \sqrt{\frac{2 \log(2) K + 2 \log(1/\delta)}{N}}, \quad (12)$$

with probability at least  $1 - \delta$ . We obtain the final bound by substituting  $K = 2\mathcal{N}(2\mathcal{X}, \ell_2, r_t(X))$  and taking a union bound on the events that inequalities (11) and (12) do not occur. □

## C Additional theoretical results

### C.1 Generalization of Lemma 1 to any element-wise activation function

**Lemma 5.** Let  $f^{(t)}$  be a  $d$ -layer NN with arbitrary activation functions at the  $t$ -th SGD iteration, demote by  $\mathbf{x}^{(t)} \in X$  the point of the training set sampled at that iteration, and set

$$\epsilon_{f^{(t)}}(\mathbf{x}, y) := \left| \frac{\partial \ell(o, y)}{\partial o} \right|_{o=f^{(t)}(\mathbf{x})}. \quad (13)$$

The Lipschitz constant of  $f^{(t)}$  at  $\mathbf{x}^{(t)}$  is

$$\frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\alpha_t \cdot \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \cdot \sigma_n(\mathbf{W}_1^{(t)}) \leq \lambda_{f^{(t)}}(\mathbf{x}^{(t)}) \leq \frac{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2}{\alpha_t \cdot \epsilon_{f^{(t)}}(\mathbf{x}^{(t)}, y^{(t)})} \cdot \sigma_1(\mathbf{W}_1^{(t)}),$$

where  $\sigma_1(\mathbf{W}_1^{(t)}) \geq \dots \geq \sigma_n(\mathbf{W}_1^{(t)}) > 0$  are the singular values of  $\mathbf{W}_1^{(t)}$ .

*Proof.* The proof proceeds almost identically with that of Lemma 1. The main difference is that the diagonal matrix  $\mathbf{S}_l^{(t)}(\mathbf{x}^{(t)})$  is redefined to yield the appropriate derivative for the activation function in question. Further, since now the function is not piece-wise linear, the bound only holds for  $\mathbf{x}^{(t)}$  (and not for the entire region  $\mathcal{R}_{\mathbf{x}^{(t)}}$  enclosing the point, as before).  $\square$

### C.2 The Lipschitz constant of the first layer

The behavior of SGD can also be indicative of the Lipschitz constant of the first layer when the training data is sufficiently diverse and the training has converged:

**Lemma 6.** Let  $f^{(t)}$  be a  $d$ -layer NN trained by SGD, let Assumption 1 hold, and further and suppose that after iteration  $\tau$ , we have

$$\frac{\|\mathbf{W}_2^{(t+1)} - \mathbf{W}_2^{(t)}\|_2}{\|\mathbf{b}_1^{(t+1)} - \mathbf{b}_1^{(t)}\|_2} + \|\mathbf{b}_1^{(t)}\|_2 \leq \vartheta \quad \text{and} \quad \|\mathbf{W}_1^{(t)} - \mathbf{W}_1^{(t')}\|_2 \leq \beta \quad \text{for all } t, t' \geq \tau.$$

Denote by  $\delta$  the minimal scalar such that, for every  $\mathbf{x} \in \mathcal{X}$ , we have  $\|\mathbf{x} - \mathbf{x}_i\|_2 \leq \delta$  for some  $\mathbf{x}_i \in X$ . Then,

$$\lambda_{f^{(t)}} \leq \frac{\vartheta + \beta}{1 - \delta}. \quad (14)$$

under the condition  $\delta < 1$ .

*Proof.* The weight matrix gradient is at a point  $\mathbf{x}$  is

$$\left( \frac{\partial f(\mathbf{x}, \mathbf{w}^{(t)})}{\partial \mathbf{W}_l^{(t)}} \right)^\top = f_{l-1}(\mathbf{x}, \mathbf{w}^{(t)}) \cdot \mathbf{W}_d^{(t)} \dots \mathbf{S}_{l+1}^{(t)}(\mathbf{x}) \mathbf{W}_{l+1}^{(t)} \mathbf{S}_l^{(t)}(\mathbf{x}).$$

Fixing

$$\left\| \left( \frac{\partial \ell(f(\mathbf{x}^{(t)}, \mathbf{w}^{(t)}), y^{(t)})}{\partial \mathbf{W}_l^{(t)}} \right)^\top \right\|_2 \left\| \left( \frac{\partial \ell(f(\mathbf{x}^{(t)}, \mathbf{w}^{(t)}), y^{(t)})}{\partial \mathbf{b}_{l-1}^{(t)}} \right)^\top \right\|_2^{-1} \leq \varrho_l(\mathbf{x}^{(t)})$$

we have that

$$\begin{aligned} \left\| \left( \frac{\partial \ell(f(\mathbf{x}^{(t)}, \mathbf{w}^{(t)}), y^{(t)})}{\partial \mathbf{W}_l^{(t)}} \right)^\top \right\| &= \|f_{l-1}(\mathbf{x}^{(t)}, \mathbf{w}^{(t)})\| \left\| \frac{\partial \ell(\hat{y}, y)}{\partial \hat{y}} \mathbf{W}_d^{(t)} \mathbf{S}_{d-1}^{(t)}(\mathbf{x}^{(t)}) \dots \mathbf{W}_{l+1}^{(t)} \mathbf{S}_l^{(t)}(\mathbf{x}^{(t)}) \right\| \\ &= \|f_{l-1}(\mathbf{x}^{(t)}, \mathbf{w}^{(t)})\|_2 \left\| \left( \frac{\partial \ell(f(\mathbf{x}^{(t)}, \mathbf{w}^{(t)}), y^{(t)})}{\partial \mathbf{b}_{l-1}^{(t)}} \right)^\top \right\|_2, \end{aligned}$$

which implies

$$\left\| f_{l-1}(\mathbf{x}^{(t)}, \mathbf{w}^{(t)}) \right\| \leq \varrho_l(\mathbf{x}^{(t)}). \quad (15)$$

Let  $\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x} \in \mathcal{S}_{n-1}} \left\| \mathbf{S}_1^{(t)}(\mathbf{x}) \mathbf{W}_1^{(t)} \mathbf{x} \right\|_2$  and fix  $\mathbf{x}^{(t')}$  to be the point in the training set that is closest to it (sampled at iteration  $t' \geq \tau$ ).

$$\lambda_{f_1^{(t)}} = \left\| \mathbf{S}_1^{(t)}(\mathbf{x}^*) \mathbf{W}_1^{(t)} \mathbf{x}^* \right\|_2 \leq \left\| \mathbf{S}_1^{(t)}(\mathbf{x}^{(t')}) \mathbf{W}_1^{(t)} \mathbf{x}^{(t')} \right\|_2 + \left\| \mathbf{S}_1^{(t)}(\mathbf{x}^*) \mathbf{W}_1^{(t)} \mathbf{x}^* - \mathbf{S}_1^{(t)}(\mathbf{x}^{(t')}) \mathbf{W}_1^{(t)} \mathbf{x}^{(t')} \right\|_2.$$

By the main assumption, we can bound the rightmost term by  $\|\mathbf{x}^* - \mathbf{x}^{(t')}\| \lambda_{f_1^{(t)}} \leq \delta \lambda_{f_1^{(t)}}$ . We thus get

$$\begin{aligned} \lambda_{f_1^{(t)}} &= \left\| \mathbf{S}_1^{(t)}(\mathbf{x}^*) \mathbf{W}_1^{(t)} \right\|_2 \leq \left\| \mathbf{S}_1^{(t)}(\mathbf{x}^{(t')}) \mathbf{W}_1^{(t)} \mathbf{x}^{(t')} \right\|_2 + \delta \lambda_{f_1^{(t)}} \\ &\leq \left\| \mathbf{S}_1^{(t)}(\mathbf{x}^{(t')}) \mathbf{W}_1^{(t')} \mathbf{x}^{(t')} \right\|_2 + \left\| \mathbf{W}_1^{(t')} - \mathbf{W}_1^{(t)} \right\|_2 + \delta \lambda_{f_1^{(t)}} \\ &\leq \left\| \mathbf{S}_1^{(t')}(\mathbf{x}^{(t')}) \mathbf{W}_1^{(t)} \mathbf{x}^{(t')} + \mathbf{b}_1^{(t')} \right\|_2 + \|\mathbf{b}_1^{(t')}\|_2 + \left\| \mathbf{W}_1^{(t')} - \mathbf{W}_1^{(t)} \right\|_2 + \delta \lambda_{f_1^{(t)}} \\ &= \varrho_2(\mathbf{x}^{(t')}) + \|\mathbf{b}_1^{(t')}\|_2 + \left\| \mathbf{W}_1^{(t')} - \mathbf{W}_1^{(t)} \right\|_2 + \delta \lambda_{f_1^{(t)}} \\ &\leq \vartheta + \left\| \mathbf{W}_1^{(t')} - \mathbf{W}_1^{(t)} \right\|_2 + \delta \lambda_{f_1^{(t)}}, \end{aligned}$$

where, due to Assumption 1,  $\mathbf{S}_1^{(t')}(\mathbf{x}^{(t')}) = \mathbf{S}_1^{(t)}(\mathbf{x}^{(t')})$ . The final bound is obtained re-arrangement and by the convergence assumption  $\left\| \mathbf{W}_1^{(t')} - \mathbf{W}_1^{(t)} \right\|_2 \leq \beta$ .  $\square$

## References

- [1] Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- [2] Huan Xu and Shie Mannor. Robustness and generalization. *Machine learning*, 86(3):391–423, 2012.
- [3] Jure Sokolić, Raja Giryes, Guillermo Sapiro, and Miguel RD Rodrigues. Robust large margin deep neural networks. *IEEE Transactions on Signal Processing*, 65(16):4265–4280, 2017.