

7. Details for the Toy Example

Suppose we train the network with n neurons for T time using gradient descent with random initialization, i.e., the network we obtain is $f_{\rho_T^n}$ using the terminology in Section 3.3. As shown by Song et al. (2018); Mei et al. (2019), $\mathcal{L}[f_{\rho_T^n}]$ is actually $\mathcal{O}(1/n + \epsilon)$ with high probability, where $\epsilon = \mathcal{L}[f_{\rho_T^\infty}]$ is the loss of the mean field limit network at training time T . Song et al. (2018) shows that $\lim_{T \rightarrow \infty} \mathcal{L}[f_{\rho_T^\infty}] = 0$ under some regularity conditions and this implies that if the training time T is sufficient, $\mathcal{L}[f_{\rho_T^\infty}]$ is generally a smaller term compared with the $\mathcal{O}(1/n)$ term.

To generate the synthesis data, we first generate a neural network $f_{\text{gen}}(\mathbf{x}) = \frac{1}{1000} \sum_{i=1}^N b_i \text{sigmoid}(\mathbf{a}_i^\top \mathbf{x})$, where \mathbf{a}_i are i.i.d. sample from a 10 dimensional standard Gaussian distribution and b_i are i.i.d. sample from a uniform distribution $\text{Unif}(-5, 5)$. The training data \mathbf{x} is also generated from a 10 dimensional standard Gaussian distribution. We choose $f_{\text{gen}}(\mathbf{x}) = y$ as the label of data. Our training data consists of 100 data points. The network we use to fit the data is $f = \frac{1}{n} \sum_{i=1}^n b'_i \tanh(\mathbf{a}'_i{}^\top \mathbf{x})$. We use network with 1000 neurons for pruning and the pruned models will not be finetuned. All networks are trained for same and sufficiently large time to converge.

8. Finding Sub-Networks on CIFAR-10/100

In this subsection, we display the results of applying our proposed algorithm to various model structures on CIFAR-10 and CIFAR-100. On CIFAR-10 and CIFAR-100, we apply our algorithm to the networks already pruned by network slimming (Liu et al., 2017) provided by Liu et al. (2019b) and show that we can further compress models which have already pruned by the L_1 regularization. We apply our algorithm on the pretrained models, and finetune the model with the same experimental setting as ImageNet.

As demonstrated in Table 4, our proposed algorithm can further compress a model pruned by Liu et al. (2019b) without or only with little drop on accuracy. For example, on the pretrained VGG19 on CIFAR-10, Liu et al. (2017) can prune 30% channels and get $93.81\% \pm 0.14\%$ accuracy. Our algorithm can prune 44% channels of the original VGG19 and get $93.78\% \pm 0.16\%$ accuracy, which is almost the same as the strong baseline number reported by Liu et al. (2019b).

DataSet	Model	Method	Prune Ratio (%)	Accuracy (%)
CIFAR10	VGG19	Liu et al. (2017)	70	93.81 ± 0.14
		Ours	56	93.78 ± 0.16
	PreResNet-164	Liu et al. (2017)	60	94.90 ± 0.04
		Ours	51	94.91 ± 0.06
		Liu et al. (2017)	40	94.71 ± 0.21
		Ours	33	94.68 ± 0.17
CIFAR100	VGG19	Liu et al. (2017)	50	73.08 ± 0.22
		Ours	44	73.05 ± 0.19
	PreResNet-164	Liu et al. (2017)	60	76.68 ± 0.35
		Ours	53	76.63 ± 0.37
		Liu et al. (2017)	40	75.73 ± 0.29
		Ours	37	75.74 ± 0.32

Table 4. Accuracy on CIFAR100 and CIFAR10. ‘‘Prune ratio’’ stands for the total percentage of channels that are pruned in the whole network. We apply our algorithm on the models pruned by Liu et al. (2017) and find that our algorithm can further prune the models. The performance of Liu et al. (2017) is reported by Liu et al. (2019b). Our reported numbers are averaged by five runs.

9. Discussion on Assumption 2 and 5

Let $\phi_j(\boldsymbol{\theta}) = \sigma(\mathbf{x}^{(j)}, \boldsymbol{\theta})/\sqrt{m}$ and $\boldsymbol{\phi}(\boldsymbol{\theta}) = [\phi_1(\boldsymbol{\theta}), \dots, \phi_m(\boldsymbol{\theta})]$ to be the vector of the outputs of the neuron $\sigma(\mathbf{x}; \boldsymbol{\theta})$ scaled by $1/\sqrt{m}$, realized on a dataset $\mathcal{D}_m := \{\mathbf{x}^{(j)}\}_{j=1}^m$. We call $\boldsymbol{\phi}(\boldsymbol{\theta})$ the feature map of $\boldsymbol{\theta}$. Given a large network $f_{[N]}(\mathbf{x}) = \sum_{i=1}^N \sigma(\mathbf{x}; \boldsymbol{\theta}_i)/N$, define the marginal polytope of the feature map to be

$$\mathcal{M}_N := \text{conv} \{ \boldsymbol{\phi}(\boldsymbol{\theta}_i) \mid i \in \{1, \dots, N\} \},$$

where conv denotes the convex hull. Then it is easy to see that Assumption 2 is equivalent to saying that $\mathbf{y} := [y^{(1)}, \dots, y^{(m)}]/\sqrt{m}$ is in the interior of the marginal polytope \mathcal{M}_N , i.e., there exists $\gamma > 0$ such that $\mathcal{B}(\mathbf{y}, \gamma) \subseteq \mathcal{M}_N$.

Here we denote by $\mathcal{B}(\boldsymbol{\mu}, r)$ the ball with radius r centered at $\boldsymbol{\mu}$. Similar to Assumption 2, Assumption 5 is equivalent to require that $\mathcal{B}(\mathbf{y}, \gamma^*) \subseteq \mathcal{M}$, where

$$\mathcal{M} := \text{conv} \{ \phi(\boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \text{supp}(\rho_T^\infty) \}.$$

We may further relax the assumption to assuming \mathbf{y} is in the relative interior (instead of interior) of \mathcal{M}_N and \mathcal{M} . However, this requires some refined analysis and we leave this as future work.

It is worth mention that when \mathcal{M} has dimension m and $f_{\rho_T^\infty}$ gives zero training loss, then assumption 5 holds. Similarly, if \mathcal{M}_N has dimension m and $f_{\rho_T^N}$ gives zero training loss, then assumption 2 holds.

10. Pruning Randomly Weighted Networks

Our theoretical analysis is also applicable for pruning randomly weighted networks. Here we give the following corollary.

Corollary 4. *Under Assumption 1 and suppose that the weights $\{\boldsymbol{\theta}_i\}$ of the large neurons $f_{[N]}(\mathbf{x})$ are i.i.d. drawn from an absolutely continuous distribution ρ_0 with a bounded support in \mathbb{R}^d , without further gradient descent training. Suppose that Assumption 5 and 6 hold for ρ_0 (changing ρ_T^∞ to ρ_0). Let S_n^{Random} be the subset obtained by the proposed greedy forward selection (2) on such $f_{[N]}$ at the n -th step. For any $\delta > 0$ and $\gamma < \gamma^*/2$, when N is sufficiently large, with probability at least $1 - \delta$, we have*

$$\mathcal{L}[f_{S_n^{\text{Random}}}] = \mathcal{O} \left(1 / (\min(1, \gamma) n)^2 \right).$$

This corollary is a special case of Theorem 3 when taking the training time to be zero ($T = 0$). And as the network is not trained, Assumption 4 are not needed for this corollary.

11. Forward Selection is Better Than Backward Elimination

A greedy backward elimination can be developed analogous to our greedy forward selection, in which we start with the full network and greedily eliminate neurons that gives the smallest increase in loss. Specifically, starting from $S_0^{\text{B}} = [N]$, we sequentially delete neurons via

$$S_{n+1}^{\text{B}} \leftarrow S_n^{\text{B}} \setminus \{i_n\}^*, \quad \text{where } i_n^* = \arg \min_{i \in S_n^{\text{B}}} \mathcal{L}[f_{S_n^{\text{B}} \setminus \{i\}}], \quad (7)$$

where \setminus denotes set minus. In this section, we demonstrate that the forward selection has significant advantages over backward elimination, from both theoretical and empirical perspectives.

Theoretical Comparison of Forward and Backward Methods Although greedy forward selection guarantees $\mathcal{O}(1/n)$ or $\mathcal{O}(1/n^2)$ error rate as we show in the paper, backward elimination does not enjoy similar theoretical guarantees. This is because the ‘‘effective search space’’ of backward elimination is more limited than that of forward selection, and gradually shrinkage over time. Specifically, at each iteration of backward elimination (7), the best neuron is chosen among S_n^{B} , which shrinks as more neurons are pruned. In contrast, the new neurons in greedy selection (2) are always selected from the full set $[N]$, which permits each neuron to be selected at every iteration, for multiple times. We now elaborate the theoretical advantages of forward selection vs. backward elimination from 1) the best achievable loss by both methods and 2) the decrease of loss across iterations.

• *On the lower bound.* In greedy forward selection, one neuron can be selected for multiple times at different iterations, while in backward elimination one neuron can only be deleted once. As a result, the best possible loss achievable by backward elimination is worse than that of greedy elimination. Specifically, because backward elimination yields a subnetwork in which each neuron appears at most once. We have an immediate lower bound of

$$\mathcal{L}[S_n^{\text{B}}] \geq \mathcal{L}_N^{\text{B}*}, \quad \forall n \in [N],$$

where

$$\mathcal{L}_N^{\text{B}*} = \min_{\boldsymbol{\alpha}} \left\{ \mathcal{L}[f_{\boldsymbol{\alpha}}] : \alpha_i = \bar{\alpha}_i / \sum_{i=1}^N \bar{\alpha}_i, \quad \bar{\alpha}_i \in \{0, 1\} \right\}.$$

In comparison, for S_n^* from forward selection (2), we have from Theorem 1 that

$$\mathcal{L}[S_n^*] = \mathcal{O}(1/n) + \mathcal{L}_N^*,$$

where \mathcal{L}_N^* equals (from Eq 3)

$$\mathcal{L}_N^* = \min_{\alpha} \left\{ \mathcal{L}[f_{\alpha}] : a_i \geq 0, \sum_{i=1}^N \alpha_i = 1 \right\}.$$

This yields a simple comparison result of

$$\mathcal{L}[S_n^B] \geq \mathcal{L}[S_n^*] + (\mathcal{L}_N^{B*} - \mathcal{L}_N^*) + \mathcal{O}(1/n).$$

Obviously, we have $\mathcal{L}_N^{B*} \geq \mathcal{L}_N^*$ because \mathcal{L}_N^* optimizes on a much larger set of α , indicating that backward elimination is inferior to forward selection. In fact, because \mathcal{L}_N^{B*} is most likely to be strictly larger than \mathcal{L}_N^* in practice, we can conclude that $\mathcal{L}[S_n^B] = \Omega(1) + \mathcal{L}_N^*$ where Ω is the Big Omega notation. This shows that it is impossible to prove bounds similar to $\mathcal{L}[S_n^*] = \mathcal{O}(1/n) + \mathcal{L}_N^*$ in Theorem 1 for backward elimination.

• *On the loss descend.* The key ingredient for proving the $\mathcal{O}(n^{-1})$ convergence of greedy forward selection is a recursive inequality that bounds $\mathcal{L}[f_{S_n}]$ at iteration n using $\mathcal{L}[f_{S_{n-1}}]$ from the previous iteration $n - 1$. Specifically, we have

$$\mathcal{L}[f_{S_n}] \leq \mathcal{L}_N^* + \frac{\mathcal{L}_N^* - \mathcal{L}[f_{S_{n-1}}]}{n} + \frac{C}{n^2}, \quad (8)$$

where $C = \max_{\mathbf{u}, \mathbf{v}} \left\{ \|\mathbf{u} - \mathbf{v}\|^2 : \mathbf{u}, \mathbf{v} \in \mathcal{M}_N \right\}$; see Appendix 12.1 for details. And inequality (8) directly implies that

$$\mathcal{L}[f_{S_n}] \leq \mathcal{L}_N^* + \frac{\mathcal{L}[f_{S_0}] - \mathcal{L}_N^*}{n}, \quad \forall n \in [N].$$

An importance reason for this inequality to hold is that the best neuron to add is selected from the whole set $[N]$ at each iteration. However, similar result does not hold for backward elimination, because the neuron to eliminate is selected from S_n^B , whose size shrinks when n grows. In fact, for backward elimination, we guarantee to find counter examples that violate a counterpart of (8), as shown in the following result, and thus fail to give the $\mathcal{O}(n^{-1})$ convergence rate.

Theorem 5. *For the S_n^B constructed by backward elimination in (7). There exists a full network $f_{[N]}(\mathbf{x}) = \sum_{i=1}^N \sigma(\mathbf{x}; \theta_i)/N$ and a dataset $\mathcal{D}_m = (\mathbf{x}^{(i)}, y^{(i)})_{i=1}^m$ that satisfies Assumption 1, 2, such that $\mathcal{L}_N^{B*} > 0$ and $\exists n \in [N]$*

$$\mathcal{L}[f_{S_{N-n}^B}] > \mathcal{L}_N^{B*} + \frac{\mathcal{L}[f_{S_N^B}] - \mathcal{L}_N^{B*}}{n},$$

In comparison, the S_n from greedy forward selection satisfies

$$\mathcal{L}[f_{S_n}] \leq \mathcal{L}_N^* + \frac{\mathcal{L}[f_{S_0}] - \mathcal{L}_N^*}{n}, \quad \forall n \in [N]. \quad (9)$$

In fact, on the same instance, we have $\mathcal{L}_N^ = 0$, and the faster rate $\mathcal{L}[f_{S_n}] \leq \mathcal{L}_N^* = \mathcal{O}(n^{-2})$ also holds for greedy forward selection.*

Proof. Suppose the data set contains 2 data points and we represent the neurons as the feature map as in section 9. Suppose that $N = 43$, $\phi(\theta_1) = [0, 1.5]$, $\phi(\theta_2) = [0, 0]$, $\phi(\theta_3) = [-0.5, 1]$, $\phi(\theta_4) = [2, 1]$ and $\phi(\theta_i) = [(-1.001)^{i-3} + 2, 1]$, $i \in \{5, 6, \dots, 43\}$ and the target $\mathbf{y} = [0, 1]$ (it is easy to construct the actual weights of neurons and data points such that the above feature maps hold). Deploying greedy backward elimination on this case gives that

$$\mathcal{L}[f_{S_{N-n}^B}] > \frac{\mathcal{L}[f_{S_N^B}] - \mathcal{L}_N^{B*}}{n} + \mathcal{L}_N^{B*},$$

for $n \in [38]$, where $\mathcal{L}_N^{B*} = \min_{n \in [N]} \mathcal{L}_{N,n}^{B*} > 0.03$. In comparison, for greedy forward selection, (9) holds from the proof of Theorem 1. In addition, on the same instance, we can verify that $\mathcal{L}_N^* = 0$, and the faster $\mathcal{O}(n^{-2})$ convergence rate also holds for greedy forward selection. In deed, the greedy forward selection is able to achieve 0 loss using two neurons (by selecting $\phi(\theta_3)$ for four times and $\phi(\theta_4)$ once). \square

Greedy Subnetwork Selection

Model	Method	Top1 Acc	FLOPs
ResNet34	Backward	73.1	2.81G
	Forward	73.5	2.64G
	Backward	72.4	2.22G
	Forward	72.9	2.07G
MobileNetV2	Backward	71.4	257M
	Forward	71.9	258M
	Backward	70.8	215M
	Forward	71.2	201M

Table 5. Comparing greedy forward selection and backward elimination on Imagenet.

Empirical Comparison of Forward and Backward Methods We compare forward selection and backward elimination to prune Resnet34 and MobilenetV2 on Imagenet. As shown in Table 5, forward selection tends to achieve better top-1 accuracy in all the cases, which is consistent with the theoretical analysis above. The experimental settings of the greedy backward elimination is the same as that of the greedy forward selection.

12. Proofs

Our proofs use the definition of the convex hulls defined in Section 9 of Appendix.

12.1. Proof of Proposition 1

The proof of Proposition 1 follows the standard argument of proving the convergence rate of Frank-Wolfe algorithm with some additional arguments. Our algorithm is not a Frank-Wolfe algorithm, but as illustrated in the subsequent proof, we can essentially use the Frank-Wolfe updates to control the error of our algorithm.

Define $\ell(\mathbf{u}) = \|\mathbf{u} - \mathbf{y}\|^2$, then the subnetwork selection problem can be viewed as solving

$$\min_{\mathbf{u} \in \mathcal{M}_N} \ell(\mathbf{u}),$$

with $\mathcal{L}_N^* = \min_{\mathbf{u} \in \mathcal{M}_N} \ell(\mathbf{u})$. And our algorithm can be viewed as starting from $\mathbf{u}^0 = 0$ and iteratively updating \mathbf{u} by

$$\mathbf{u}^k = (1 - \xi_k)\mathbf{u}^{k-1} + \xi_k \mathbf{q}^k, \quad \mathbf{q}^k = \arg \min_{\mathbf{q} \in \text{Vert}(\mathcal{M}_N)} \ell((1 - \xi_k)\mathbf{u}^{k-1} + \xi_k \mathbf{q}), \quad (10)$$

where $\text{Vert}(\mathcal{M}_N) := \{\phi(\theta_1), \dots, \phi(\theta_N)\}$ denotes the vertices of \mathcal{M}_N , and we shall take $\xi_k = 1/k$. We aim to prove that $\ell(\mathbf{u}^k) = O(1/k) + \mathcal{L}_N^*$. Our proof can be easily extended to general convex functions $\ell(\cdot)$ and different ξ_k schemes.

By the convexity and the quadratic form of $\ell(\cdot)$, for any \mathbf{s} , we have

$$\ell(\mathbf{s}) \geq \ell(\mathbf{u}^{k-1}) + \nabla \ell(\mathbf{u}^{k-1})^\top (\mathbf{s} - \mathbf{u}^{k-1}) \quad (11)$$

$$\ell(\mathbf{s}) \leq \ell(\mathbf{u}^{k-1}) + \nabla \ell(\mathbf{u}^{k-1})^\top (\mathbf{s} - \mathbf{u}^{k-1}) + \|\mathbf{s} - \mathbf{u}^{k-1}\|^2. \quad (12)$$

Minimizing \mathbf{s} in \mathcal{M}_N on both sides of (11), we have

$$\begin{aligned} \mathcal{L}_N^* &= \min_{\mathbf{s} \in \mathcal{M}_N} \ell(\mathbf{s}) \geq \min_{\mathbf{s} \in \mathcal{M}_N} \{\ell(\mathbf{u}^{k-1}) + \nabla \ell(\mathbf{u}^{k-1})^\top (\mathbf{s} - \mathbf{u}^{k-1})\} \\ &= \ell(\mathbf{u}^{k-1}) + \nabla \ell(\mathbf{u}^{k-1})^\top (\mathbf{s}^k - \mathbf{u}^{k-1}). \end{aligned} \quad (13)$$

Here we define

$$\begin{aligned} \mathbf{s}^k &= \arg \min_{\mathbf{s} \in \mathcal{M}_N} \nabla \ell(\mathbf{u}^{k-1})^\top (\mathbf{s} - \mathbf{u}^{k-1}) \\ &= \arg \min_{\mathbf{s} \in \text{Vert}(\mathcal{M}_N)} \nabla \ell(\mathbf{u}^{k-1})^\top (\mathbf{s} - \mathbf{u}^{k-1}), \end{aligned} \quad (14)$$

where the second equation holds because we optimize a linear objective on a convex polytope \mathcal{M}_N and hence the solution must be achieved on the vertices $\text{Vert}(\mathcal{M}_N)$. Note that if we update \mathbf{u}^k by $\mathbf{u}^k = (1 - \xi_k)\mathbf{u}^{k-1} + \xi_k \mathbf{s}^k$, we would get the standard Frank-Wolfe (or conditional gradient) algorithm. The difference between our method and Frank-Wolfe is that we greedily minimize the loss $\ell(\mathbf{u}^k)$, while the Frank-Wolfe minimizes the linear approximation in (14).

Define $D_{\mathcal{M}_N} := \max_{\mathbf{u}, \mathbf{v}} \{\|\mathbf{u} - \mathbf{v}\| : \mathbf{u}, \mathbf{v} \in \mathcal{M}_N\}$ to be the diameter of \mathcal{M}_N . Following (17), we have

$$\begin{aligned} \ell(\mathbf{u}^k) &= \min_{\mathbf{q} \in \text{Vert}(\mathcal{M}_N)} \ell((1 - \xi_k)\mathbf{u}^{k-1} + \xi_k \mathbf{q}) \\ &\leq \ell((1 - \xi_k)\mathbf{u}^{k-1} + \xi_k \mathbf{s}^k) \\ &\leq \ell(\mathbf{u}^{k-1}) + \xi_k \nabla \ell(\mathbf{u}^{k-1})^\top (\mathbf{s}^k - \mathbf{u}^{k-1}) + C\xi_k^2 \end{aligned} \quad (15)$$

$$\leq (1 - \xi_k)\ell(\mathbf{u}^{k-1}) + \xi_k \mathcal{L}_N^* + C\xi_k^2, \quad (16)$$

where we define $C := D_{\mathcal{M}_N}^2$, (15) follows (12), and (16) follows (13). Rearranging this, we get

$$\ell(\mathbf{u}^k) - \mathcal{L}_N^* - C\xi_k \leq (1 - \xi_k) (\ell(\mathbf{u}^{k-1}) - \mathcal{L}_N^* - C\xi_k)$$

By iteratively applying the above inequality, we have

$$\ell(\mathbf{u}^k) - \mathcal{L}_N^* - C\xi_k \leq \left(\prod_{i=1}^k (1 - \xi_i) \right) (\ell(\mathbf{u}^0) - \mathcal{L}_N^* - C\xi_1).$$

Taking $\xi_k = 1/k$. We get

$$\ell(\mathbf{u}^k) - \mathcal{L}_N^* - \frac{C}{k} \leq \frac{1}{k} (\ell(\mathbf{u}^0) - \mathcal{L}_N^* - C).$$

And thus

$$\ell(\mathbf{u}^k) \leq \frac{1}{k} (\ell(\mathbf{u}^0) - \mathcal{L}_N^*) + \mathcal{L}_N^* = \mathcal{O}\left(\frac{1}{k}\right) + \mathcal{L}_N^*.$$

This completes the proof.

12.2. Proof of Theorem 2

The proof leverages the idea from the proof of Proposition 1 of [Chen et al. \(2012\)](#) for analyzing their *Herding* algorithm, but contains some extra nontrivial argument.

Following the proof of Proposition 1, our problem can be viewed as

$$\min_{\mathbf{u} \in \mathcal{M}_N} \left\{ \ell(\mathbf{u}) := \|\mathbf{u} - \mathbf{y}\|^2 \right\},$$

with $\mathcal{L}_N^* = \min_{\mathbf{u} \in \mathcal{M}_N} \ell(\mathbf{u})$, our greedy algorithm can be viewed as starting from $\mathbf{u}^0 = 0$ and iteratively updating \mathbf{u} by

$$\mathbf{u}^k = \frac{k-1}{k} \mathbf{u}^{k-1} + \frac{1}{k} \mathbf{q}^k, \quad \mathbf{q}^k = \arg \min_{\mathbf{q} \in \text{Vert}(\mathcal{M}_N)} \left\| \frac{k-1}{k} \mathbf{u}^{k-1} + \frac{1}{k} \mathbf{q} - \mathbf{y} \right\|^2 \quad (17)$$

where $\text{Vert}(\mathcal{M}_N) := \{\phi(\boldsymbol{\theta}_1), \dots, \phi(\boldsymbol{\theta}_N)\}$ denotes the vertices of \mathcal{M}_N . We aim to prove that $\ell(\mathbf{u}^k) = O(1/(k \max(1, \gamma))^2)$, under Assumption 2.

Define $\mathbf{w}^k = k(\mathbf{y} - \mathbf{u}^k)$, then $\ell(\mathbf{u}^k) = \|\mathbf{w}^k\|^2 / k^2$. Therefore, it is sufficient to prove that $\|\mathbf{w}^k\| = \mathcal{O}(1/(\max(1, \gamma)))$.

Similar to the proof of Proposition 1, we define

$$\begin{aligned}
 \mathbf{s}^{k+1} &= \arg \min_{\mathbf{s} \in \mathcal{M}_N} \nabla \ell(\mathbf{u}^k)^\top (\mathbf{s} - \mathbf{u}^k) \\
 &= \arg \min_{\mathbf{s} \in \mathcal{M}_N} \nabla \ell(\mathbf{u}^k)^\top \mathbf{s} \\
 &= \arg \min_{\mathbf{s} \in \mathcal{M}_N} \langle \mathbf{w}^k, \mathbf{s} \rangle. \\
 &= \arg \min_{\mathbf{s} \in \mathcal{M}_N} \langle \mathbf{w}^k, (\mathbf{s} - \mathbf{y}) \rangle.
 \end{aligned}$$

Because $\mathcal{B}(\mathbf{y}, \gamma)$ is included in \mathcal{M}_N by Assumption 2, we have $\mathbf{s}' := \mathbf{y} - \gamma \mathbf{w}^k / \|\mathbf{w}^k\| \in \mathcal{M}_N$. Therefore

$$\langle \mathbf{w}^k, (\mathbf{s}^{k+1} - \mathbf{y}) \rangle = \min_{\mathbf{s} \in \mathcal{M}_N} \langle \mathbf{w}^k, (\mathbf{s} - \mathbf{y}) \rangle \leq \langle \mathbf{w}^k, (\mathbf{s}' - \mathbf{y}) \rangle = -\gamma \|\mathbf{w}^k\|.$$

Note that

$$\begin{aligned}
 \|\mathbf{w}^{k+1}\|^2 &= \min_{\mathbf{q} \in \text{Vert}(\mathcal{M}_N)} \|k\mathbf{u}^k + \mathbf{q} - (k+1)\mathbf{y}\|^2 \\
 &= \min_{\mathbf{q} \in \text{Vert}(\mathcal{M}_N)} \|\mathbf{w}^k + \mathbf{q} - \mathbf{y}\|^2 \\
 &\leq \|\mathbf{w}^k + \mathbf{s}^{k+1} - \mathbf{y}\|^2 \\
 &= \|\mathbf{w}^k\|^2 + 2\langle \mathbf{w}^k, (\mathbf{s}^{k+1} - \mathbf{y}) \rangle + \|\mathbf{s}^{k+1} - \mathbf{y}\|^2 \\
 &\leq \|\mathbf{w}^k\|^2 - 2\gamma \|\mathbf{w}^k\| + D_{\mathcal{M}_N}^2,
 \end{aligned}$$

where $D_{\mathcal{M}_N}$ is the diameter of \mathcal{M}_N . Because $\mathbf{w}^0 = 0$, using Lemma 6, we have

$$\|\mathbf{w}^k\| \leq \max(D_{\mathcal{M}_N}, D_{\mathcal{M}_N}^2/2, D_{\mathcal{M}_N}^2/(2\gamma)) = \mathcal{O}\left(\frac{1}{\min(1, \gamma)}\right), \quad \forall k = 1, 2, \dots,$$

This proves that $\ell(\mathbf{u}^k) = \frac{\|\mathbf{w}^k\|^2}{k^2} = \mathcal{O}\left(\frac{1}{k^2 \min(1, \gamma)^2}\right)$.

Lemma 6. Assume $\{z_k\}_{k \geq 0}$ is a sequence of numbers satisfying $z_0 = 0$ and

$$|z_{k+1}|^2 \leq |z_k|^2 - 2\gamma|z_k| + C, \quad \forall k = 0, 1, 2, \dots$$

where C and γ are two positive numbers. Then we have $|z_k| \leq \max(\sqrt{C}, C/2, C/(2\gamma))$ for all $k = 0, 1, 2, \dots$

Proof. We prove $|z_k| \leq \max(\sqrt{C}, C/2, C/(2\gamma)) := u_*$ by induction on k . Because $z_0 = 0$, the result holds for $k = 0$. Assume $|z_k| \leq u_*$, we want to prove that $|z_{k+1}| \leq u_*$ also holds.

Define $f(z) = z^2 - 2\gamma z + C$. Note that the maximum of $f(z)$ on an interval is always achieved on the vertices, because $f(z)$ is convex.

Case 1: If $|z_k| \leq C/(2\gamma)$, then we have

$$|z_{k+1}|^2 \leq f(|z_k|) \leq \max_z \left\{ f(z) : z \in [0, C/(2\gamma)] \right\} = \max \left\{ f(0), f(C/(2\gamma)) \right\} = \max \left\{ C, C^2/(4\gamma^2) \right\} \leq u_*^2.$$

Case 2: If $|z_k| \geq C/(2\gamma)$, then we have

$$|z_{k+1}|^2 \leq |z_k|^2 - 2\gamma|z_k| + C \leq |z_k|^2 \leq u_*^2.$$

In both cases, we have $|z_{k+1}| \leq u_*$. This completes the proof. □

12.3. Proof of Theorem 3

We first introduce the following Lemmas.

Lemma 7. *Under the Assumption 1, 3, 4, 5 and 6. For any $\delta > 0$, when N is sufficient large, with probability at least $1 - \delta$,*

$$\mathcal{B}\left(\mathbf{y}, \frac{1}{2}\gamma^*\right) \subseteq \text{conv}\{\phi(\boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \text{supp}(\rho_T^N)\}.$$

Here ρ_T^N is the distribution of the weight of the large network with N neurons trained by gradient descent.

12.3.1. PROOF OF THEOREM 3

The above lemmas directly imply Theorem 3.

12.3.2. PROOF OF LEMMA 7

In this proof, we simplify the statement that ‘for any $\delta > 0$, when N is sufficiently large, event E holds with probability at least $1 - \delta$ ’ by ‘when N is sufficiently large, with high probability, event E holds’.

By the Assumption 5, there exists $\gamma^* > 0$ such that

$$\mathcal{B}(\mathbf{y}, \gamma^*) \subseteq \text{conv}\{\phi(\boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \text{supp}(\rho_T^\infty)\} = \mathcal{M}.$$

Given any $\boldsymbol{\theta} \in \text{supp}(\rho_T^\infty)$, define

$$\phi^N(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}' \in \text{supp}(\rho_T^N)} \|\phi(\boldsymbol{\theta}') - \phi(\boldsymbol{\theta})\|$$

where $\phi^N(\boldsymbol{\theta})$ is the best approximation of $\phi(\boldsymbol{\theta})$ using the points $\phi(\boldsymbol{\theta}_i)$, $\boldsymbol{\theta}_i \in \text{supp}(\rho_T^N)$.

Using Lemma 11, by choosing $\epsilon = \gamma^*/6$, when N is sufficiently large, we have

$$\sup_{\boldsymbol{\theta} \in \text{supp}(\rho_T^\infty)} \|\phi(\boldsymbol{\theta}) - \phi^N(\boldsymbol{\theta})\| \leq \gamma^*/6, \quad (18)$$

with high probability. (18) implies that \mathcal{M}_N can approximate \mathcal{M} for large N . Since \mathcal{M} is assumed to contain the ball centered at \mathbf{y} with radius γ^* , as \mathcal{M}_N approximates \mathcal{M} , intuitively \mathcal{M}_N would also contain the ball centered at \mathbf{y} with a smaller radius. And below we give a rigorous proof for this intuition.

Step 1: $\|\hat{\mathbf{y}} - \mathbf{y}\| \leq \gamma^*/6$. When N is sufficiently large, with high probability, we have

$$\|\hat{\mathbf{y}} - \mathbf{y}\| \leq \sum_{i=1}^M q_i \|\phi^N(\boldsymbol{\theta}_i^*) - \phi(\boldsymbol{\theta}_i^*)\| \leq \gamma^*/6.$$

Step 2 $\mathcal{B}(\hat{\mathbf{y}}, \frac{5}{6}\gamma^*) \subseteq \mathcal{M}$ By step one, with high probability, $\|\hat{\mathbf{y}} - \mathbf{y}\| \leq \gamma^*/4$, which implies that $\hat{\mathbf{y}} \in \mathcal{B}(\mathbf{y}, \gamma^*/4) \subseteq \mathcal{B}(\mathbf{y}, \gamma^*) \subseteq \mathcal{M}$. Also, for any $A \in \partial\mathcal{M}$ (here $\partial\mathcal{M}$ denotes the boundary of \mathcal{M}), we have

$$\|\hat{\mathbf{y}} - A\| \geq \|\mathbf{y} - A\| - \|\mathbf{y} - \hat{\mathbf{y}}\| \geq \gamma^* - \gamma^*/4.$$

This gives that $\mathcal{B}(\hat{\mathbf{y}}, \frac{5}{6}\gamma^*) \subseteq \mathcal{M}$.

Step 3 $\mathcal{B}(\hat{\mathbf{y}}, \frac{2}{3}\gamma^*) \subseteq \mathcal{M}_N$ Notice that $\hat{\mathbf{y}}$ is a point in \mathbb{R}^m and suppose that A belongs to the boundary of \mathcal{M}_N (denoted by $\partial\mathcal{M}_N$) such that

$$\|\hat{\mathbf{y}} - A\| = \min_{\tilde{A} \in \partial\mathcal{M}_N} \|\hat{\mathbf{y}} - \tilde{A}\|.$$

We prove by contradiction. Suppose that we have $\|\hat{\mathbf{y}} - A\| < \frac{2}{3}\gamma^*$.

Using support hyperplane theorem, there exists a hyperplane $P = \{\mathbf{u} : \langle \mathbf{u} - A, \mathbf{v} \rangle = 0\}$ for some nonempty vector \mathbf{v} , such that $A \in P$ and

$$\sup_{\mathbf{q} \in \mathcal{M}_N} \langle \mathbf{q}, \mathbf{v} \rangle \leq \langle A, \mathbf{v} \rangle.$$

We choose $A' \in P$ such that $A' - \hat{\mathbf{y}} \perp P$ (A and A' can be the same point). Notice that

$$\|\hat{\mathbf{y}} - A'\|^2 = \|\hat{\mathbf{y}} - A + A - A'\|^2 = \|\hat{\mathbf{y}} - A\|^2 + \|A - A'\|^2 + 2\langle \hat{\mathbf{y}} - A, A - A' \rangle.$$

Since $A' - \hat{\mathbf{y}} \perp P$ and $A, A' \in P$, we have $\langle \hat{\mathbf{y}} - A, A - A' \rangle = 0$ and thus $\|\hat{\mathbf{y}} - A'\| \leq \|\hat{\mathbf{y}} - A\| < \frac{2}{3}\gamma^*$. We have

$$A' \in \mathcal{B}(\hat{\mathbf{y}}, \|\hat{\mathbf{y}} - A\|) \subseteq \mathcal{B}\left(\hat{\mathbf{y}}, \frac{2}{3}\gamma^*\right) \subseteq \mathcal{B}\left(\hat{\mathbf{y}}, \frac{5}{6}\gamma^*\right) \subseteq \mathcal{M}.$$

Notice that as both $\hat{\mathbf{y}}, A' \in \mathcal{M}$ we choose $\lambda \geq 1$ such that $\hat{\mathbf{y}} + \lambda(A' - \hat{\mathbf{y}}) \in \partial\mathcal{M}$, where $\partial\mathcal{M}$ denotes the boundary of \mathcal{M} . Define $B = \hat{\mathbf{y}} + \lambda(A' - \hat{\mathbf{y}})$. As we have shown that $\mathcal{B}\left(\hat{\mathbf{y}}, \frac{5}{6}\gamma^*\right) \subseteq \mathcal{M}$, we have $\|\hat{\mathbf{y}} - B\| \geq \frac{5}{6}\gamma^*$. And thus

$$\begin{aligned} \|B - A'\| &= \|B - \hat{\mathbf{y}}\| - \|\hat{\mathbf{y}} - A'\| \\ &> \frac{5}{6}\gamma^* - \frac{2}{3}\gamma^* \\ &> \frac{1}{6}\gamma^*. \end{aligned}$$

Also notice that

$$\begin{aligned} \langle B - A, \mathbf{v} \rangle &= \langle \hat{\mathbf{y}} + \lambda(A' - \hat{\mathbf{y}}) - A, \mathbf{v} \rangle \\ &= (1 - \lambda)\langle \hat{\mathbf{y}} - A, \mathbf{v} \rangle + \lambda\langle A' - A, \mathbf{v} \rangle \\ &= (1 - \lambda)\langle \hat{\mathbf{y}} - A, \mathbf{v} \rangle \\ &\geq 0. \end{aligned}$$

This implies that B and \mathcal{M} are on different side of P .

With high probability, we are able to find $D \in \{\phi(\boldsymbol{\theta}); \boldsymbol{\theta} \in \text{supp}(\rho_T^N)\}$ such that

$$\|D - B\| \leq \frac{\gamma^*}{6}.$$

By the definition, $D \in \mathcal{M}_N$ and thus $\langle D - A, \mathbf{v} \rangle \leq 0$ as shown by the supporting hyperplane theorem. Also remind that $\langle B - A, \mathbf{v} \rangle \geq 0$. These allow us to choose $\lambda' \in [0, 1]$ such that

$$\langle \lambda'D + (1 - \lambda')B - A, \mathbf{v} \rangle = 0.$$

We define $E = \lambda'D + (1 - \lambda')B$ and thus $E \in P$. Notice that

$$\|B - E\| = \|B - \lambda'D - (1 - \lambda')B\| = \lambda'\|B - D\| \leq \|B - D\| \leq \frac{\gamma^*}{6}.$$

Also,

$$\|B - E\|^2 = \|B - A' + A' - E\|^2 = \|B - A'\|^2 + \|A' - E\|^2 + 2\langle B - A', A' - E \rangle.$$

As $B - A' \perp P$ and $A', E \in P$, we have $\langle B - A', A' - E \rangle = 0$, which implies that $\|B - E\| \geq \|B - A'\| > \frac{1}{6}\gamma^*$, which makes contradiction.

Step 4 $\mathcal{B}(\mathbf{y}, \frac{1}{2}\gamma^*) \subseteq \mathcal{M}_N$ As for sufficiently large N , we have $\|\hat{\mathbf{y}} - \mathbf{y}\| \leq \frac{1}{6}\gamma^*$ and thus

$$\mathcal{B}\left(\mathbf{y}, \frac{1}{2}\gamma^*\right) \subseteq \mathcal{B}\left(\hat{\mathbf{y}}, \frac{2}{3}\gamma^*\right) \subseteq \mathcal{M}_N.$$

13. Technical Lemmas

Lemma 8. Under assumption 1 and 3, for any N , at training time $T < \infty$, for any $\boldsymbol{\theta} \in \text{supp}(\rho_T^N)$ or $\boldsymbol{\theta} \in \text{supp}(\rho_T^\infty)$, we have $\|\boldsymbol{\theta}\| \leq C$, $\|\phi(\boldsymbol{\theta})\| \leq C$ and $\|\phi(\boldsymbol{\theta})\|_{Lip} \leq C$ for some constant $C < \infty$.

Lemma 9. Suppose $\theta_i \in \mathbb{R}^d$, $i = 1, \dots, N$ are i.i.d. samples from some distribution ρ and $\Omega \subseteq \mathbb{R}^d$ is bounded. For any radius $r_B > 0$ and $\delta > 0$, define the following two sets

$$A = \left\{ \theta_B \in \Omega \mid \mathbb{P}_{\theta \sim \rho} (\theta \in \mathcal{B}(\theta_B, r_B)) > \frac{4}{N} ((d+1) \log(2N) + \log(8/\delta)) \right\}$$

$$B = \left\{ \theta_B \in \Omega \mid \left\| \theta_B - \theta_B^N \right\| \leq r_B \right\},$$

where $\theta_B^N = \arg \min_{\theta' \in \{\theta_i\}_{i=1}^N} \left\| \theta_B - \theta' \right\|$. With probability at least $1 - \delta$, $A \subseteq B$.

Lemma 10. For any $\delta > 0$ and $\epsilon > 0$, when N is sufficiently large (N depends on δ), with probability at least $1 - \delta$, we have

$$\sup_{\theta \in \text{supp}(\rho_T^\infty)} \left\| \phi(\theta) - \bar{\phi}^N(\theta) \right\| \leq \epsilon,$$

where $\bar{\phi}^N(\theta) = \arg \min_{\phi(\bar{\theta}') \in \{\phi(\bar{\theta}_i)\}_{i=1}^N} \left\| \phi(\bar{\theta}') - \phi(\theta) \right\|$ and $\bar{\theta}_i$ are i.i.d. samples from ρ_T^∞ .

Lemma 11. For any $\delta > 0$ and $\epsilon > 0$, when N is sufficiently large (N depends on δ), with probability at least $1 - \delta$, we have

$$\sup_{\theta \in \text{supp}(\rho_T^\infty)} \left\| \phi(\theta) - \phi^N(\theta) \right\| \leq \epsilon,$$

where $\phi^N(\theta) = \arg \min_{\theta' \in \text{supp}(\rho_T^N)} \left\| \phi(\theta') - \phi(\theta) \right\|$.

13.1. Proof of Lemma 8

We prove the case of training network with N neurons. Notice that

$$\begin{aligned} \left\| \frac{\partial}{\partial t} \theta(t) \right\| &= \left\| \mathbf{g}[\theta(t), \rho_t^N] \right\| \\ &= \left\| \mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}} (y - f_{\rho_t^N}(\mathbf{x})) \nabla_{\theta} \sigma(\theta(t), \mathbf{x}) \right\| \\ &\leq \sqrt{\mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}} (y - f_{\rho_t^N}(\mathbf{x}))^2} \sqrt{\mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}} \|\nabla_{\theta} \sigma(\theta(t), \mathbf{x})\|^2} \\ &\leq \sqrt{\mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}} (y - f_{\rho_0^N}(\mathbf{x}))^2} \sqrt{\mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}} \|\nabla_{\theta} \sigma(\theta(t), \mathbf{x})\|^2} \end{aligned}$$

Notice that by the assumption 1, we have $\sqrt{\mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}} (y - f_{\rho_0^N}(\mathbf{x}))^2} \leq C$. Remind that $\theta(t) = [\mathbf{a}(t), b(t)]$, $\sigma(\theta(t), \mathbf{x}) = b(t)\sigma_+(\mathbf{a}^\top(t)\mathbf{x})$. Thus we have

$$\left| \frac{\partial}{\partial t} b(t) \right| \leq C \|\sigma_+\|_\infty.$$

And thus for any $i \in \{1, \dots, N\}$, $\sup_{t \in [0, T]} \|b_i(t)\| \leq \int_0^T \left\| \frac{\partial}{\partial t} b_i(s) \right\| ds \leq TC$. Also

$$\begin{aligned} \left\| \frac{\partial}{\partial t} \mathbf{a}(t) \right\| &\leq C|b(t)| \|\sigma'_+\|_\infty \sqrt{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \|\mathbf{x}\|^2} \\ &\leq TC. \end{aligned}$$

By assumption 3, that $\|\theta_0(t)\| \leq C$, we have

$$\sup_{t \in [0, T]} \|\theta_i(t)\| \leq \int_0^T \left\| \frac{\partial}{\partial t} \theta_i(s) \right\| ds \leq T^2 C.$$

Notice that this also holds to training the network with infinite number of neurons. Notice that $\|\phi(\boldsymbol{\theta})\| = \sqrt{\frac{1}{m} \sum_{j=1}^m \sigma^2(\boldsymbol{\theta}, \mathbf{x}^{(j)})} \leq CT$. And

$$\begin{aligned} \|\phi(\boldsymbol{\theta})\|_{\text{Lip}} &= \sup_{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2} \frac{\|\phi(\boldsymbol{\theta}_1) - \phi(\boldsymbol{\theta}_2)\|}{\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|} \\ &= \sup_{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2} \frac{\sqrt{\frac{1}{m} \sum_{j=1}^m (\sigma(\boldsymbol{\theta}_1, \mathbf{x}^{(j)}) - \sigma(\boldsymbol{\theta}_2, \mathbf{x}^{(j)}))^2}}{\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|} \\ &\leq TC \|\sigma_+\|_{\text{Lip}} + \|\sigma_+\|_{\infty}. \end{aligned}$$

Thus given any $T < \infty$, all those three quantities can be bounded by some constant.

13.2. Proof of Lemma 9

The following proof follows line 1 and line 2 of the proof of Lemma 16 of (Chaudhuri & Dasgupta, 2010).

Define $g_{\boldsymbol{\theta}_B}(\boldsymbol{\theta}) = \mathbb{I}\{\boldsymbol{\theta} \in \mathcal{B}(\boldsymbol{\theta}_B, r_B)\}$ and $\beta_N = \sqrt{(4/N)(d_{\text{VC}} \log 2N + \log(8/\delta))}$, where d_{VC} is the VC dimension of the function class $\mathcal{G} = \{g_{\boldsymbol{\theta}_B}, \boldsymbol{\theta}_B \in \Omega\}$ and thus $d_{\text{VC}} \leq d + 1$ (Dudley, 1979). Let $\mathbb{E}g_{\boldsymbol{\theta}_B} = \mathbb{P}_{\boldsymbol{\theta} \sim \rho}(\boldsymbol{\theta} \in \mathcal{B}(\boldsymbol{\theta}_B, r_B))$ and $\mathbb{E}_N g_{\boldsymbol{\theta}_B} = \sum_{i=1}^N g_{\boldsymbol{\theta}_B}(\boldsymbol{\theta}_i)/N$. So

$$A = \{\boldsymbol{\theta}_B \mid \mathbb{E}g_{\boldsymbol{\theta}_B} > \beta_N^2\}$$

and we further define

$$A_2 = \{\boldsymbol{\theta}_B \mid \mathbb{E}_N g_{\boldsymbol{\theta}_B} > 0\}.$$

From theorem 15 of (Chaudhuri & Dasgupta, 2010) (which is a rephrase of the generalization bound), we know that: for any $\delta > 0$, with probability at least $1 - \delta$, the following holds for all $g_{\boldsymbol{\theta}_B} \in \mathcal{G}$,

$$\mathbb{E}g_{\boldsymbol{\theta}_B} - \mathbb{E}_N g_{\boldsymbol{\theta}_B} \leq \beta_N \sqrt{\mathbb{E}g_{\boldsymbol{\theta}_B}} \quad (19)$$

Notice that for any $g_{\boldsymbol{\theta}_B}$ which satisfies (19),

$$\mathbb{E}g_{\boldsymbol{\theta}_B} > \beta_N^2 \Rightarrow \mathbb{E}_N g_{\boldsymbol{\theta}_B} > 0$$

So this means: for any $\delta > 0$, with probability at least $1 - \delta$,

$$A \subseteq A_2 = B$$

where the last equality follows from the following:

$$A_2 = \{\boldsymbol{\theta}_B \mid \mathbb{E}_N g_{\boldsymbol{\theta}_B} > 0\} = \{\text{there exists some } \boldsymbol{\theta}_i \text{ such that } \boldsymbol{\theta}_i \in \mathcal{B}(\boldsymbol{\theta}_B, r_B)\} = B$$

13.3. Proof of Lemma 10

Given $\epsilon > 0$, we choose r_0 sufficiently small such that $Cr_0 \leq \epsilon$ (here C is some constant defined in Lemma 8). For this choice of r_0 , given the corresponding p_0 (defined in assumption 6), for any $\delta > 0$, there exists $N(\delta)$ such that $\forall N \geq N(\delta)$, we have

$$p_0 > \frac{4}{N} ((d+1) \log(2N) + \log(8/\delta)) := \beta_N^2.$$

And thus from assumption 6, we have

$$\forall \boldsymbol{\theta} \in \text{supp}(\rho_T^\infty), \mathbb{P}_{\boldsymbol{\theta}' \sim \rho_T^\infty}(\boldsymbol{\theta}' \in \mathcal{B}(\boldsymbol{\theta}, r_0)) \geq p_0 > \beta_N^2.$$

This implies

$$\text{supp}(\rho_T^\infty) \subseteq A = \{\boldsymbol{\theta}_B \in \Omega \mid \mathbb{P}_{\boldsymbol{\theta} \sim \rho}(\boldsymbol{\theta} \in \mathcal{B}(\boldsymbol{\theta}_B, r_0)) > \beta_N^2\}$$

From Lemma 9 (set $r_B = r_0$), we know: with probability at least $1 - \delta$,

$$A \subseteq B = \left\{ \boldsymbol{\theta}_B \in \Omega \mid \left\| \boldsymbol{\theta}_B - \boldsymbol{\theta}_B^N \right\| \leq r_0 \right\},$$

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

$$\text{supp}(\rho_T^\infty) \subseteq B$$

and this means: with probability at least $1 - \delta$, we have

$$\forall \boldsymbol{\theta} \in \text{supp}(\rho_T^\infty), \quad \|\boldsymbol{\theta} - \boldsymbol{\theta}^N\| \leq r_0.$$

The result concludes from

$$\begin{aligned} & \sup_{\boldsymbol{\theta} \in \text{supp}(\rho_T^\infty)} \|\boldsymbol{\phi}(\boldsymbol{\theta}) - \boldsymbol{\phi}^N(\boldsymbol{\theta})\| \\ & \leq \sup_{\boldsymbol{\theta} \in \text{supp}(\rho_T^\infty)} \|\boldsymbol{\phi}(\boldsymbol{\theta}) - \boldsymbol{\phi}(\boldsymbol{\theta}^N)\| \\ & \leq \sup_{\boldsymbol{\theta} \in \text{supp}(\rho_T^\infty)} C \|\boldsymbol{\theta} - \boldsymbol{\theta}^N\| \\ & \leq Cr_0 \leq \epsilon. \end{aligned}$$

Here the last inequality uses Lemma 8.

13.4. Proof of Lemma 11

In this proof, we simplify the statement that ‘for any $\delta > 0$, when N is sufficiently large, event E holds with probability at least $1 - \delta$ ’ by ‘when N is sufficiently large, with high probability, event E holds’.

Suppose that $\boldsymbol{\theta}_i, i \in [N]$ is the weight of neurons of network $f_{\rho_T^\infty}$. Given any $\boldsymbol{\theta} \in \text{supp}(\rho_T^\infty)$, define

$$\boldsymbol{\phi}^N(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\phi}(\boldsymbol{\theta}') \in \text{Vert}(\mathcal{M}_N)} \|\boldsymbol{\phi}(\boldsymbol{\theta}') - \boldsymbol{\phi}(\boldsymbol{\theta})\|.$$

Notice that the training dynamics of the network with N neurons can be characterized by

$$\begin{aligned} \frac{\partial}{\partial t} \boldsymbol{\theta}_i(t) &= \mathbf{g}[\boldsymbol{\theta}_i(t), \rho_i^N], \\ \boldsymbol{\theta}_i(0) &\stackrel{\text{i.i.d.}}{\sim} \rho_0. \end{aligned}$$

Here $\mathbf{g}[\boldsymbol{\theta}, \rho] = \mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}} (y - f_\rho(\mathbf{x})) \nabla_{\boldsymbol{\theta}} \sigma(\boldsymbol{\theta}, \mathbf{x})$. We define the following coupling dynamics:

$$\begin{aligned} \frac{\partial}{\partial t} \bar{\boldsymbol{\theta}}_i(t) &= \mathbf{g}[\bar{\boldsymbol{\theta}}_i(t), \rho_i^\infty], \\ \bar{\boldsymbol{\theta}}_i(0) &= \boldsymbol{\theta}_i(0). \end{aligned}$$

Notice that at any time t , $\bar{\boldsymbol{\theta}}_i(t)$ can be viewed as i.i.d. sample from ρ_i^∞ . We define $\hat{\rho}_i^N(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \delta_{\bar{\boldsymbol{\theta}}_i(t)}(\boldsymbol{\theta})$. Notice that by our definition $\boldsymbol{\theta}_i = \boldsymbol{\theta}_i(T)$ and we also define $\bar{\boldsymbol{\theta}}_i = \bar{\boldsymbol{\theta}}_i(T)$. Using the propagation of chaos argument as [Mei et al. \(2019\)](#) (Proposition 2 of Appendix B.2), for any $T < \infty$, for any $\delta > 0$, we have

$$\sup_{t \in [0, T]} \max_{i \in \{1, \dots, N\}} \|\bar{\boldsymbol{\theta}}_i(t) - \boldsymbol{\theta}_i(t)\| \leq \frac{C}{\sqrt{N}} \left(\sqrt{\log N} + \sqrt{\log 1/\delta} \right).$$

By Lemma 10 and the bound above, when N is sufficiently large, with high probability, we have

$$\begin{aligned} \sup_{\boldsymbol{\theta} \in \text{supp}(\rho_T^\infty)} \|\boldsymbol{\phi}(\boldsymbol{\theta}) - \bar{\boldsymbol{\phi}}^N(\boldsymbol{\theta})\| &\leq \epsilon/2 \\ \max_{i \in [N]} \|\bar{\boldsymbol{\theta}}_i(T) - \boldsymbol{\theta}_i(T)\| &\leq \frac{\epsilon}{2C}, \end{aligned}$$

where $C = \|\boldsymbol{\phi}\|_{\text{Lip}}$ and

$$\bar{\boldsymbol{\phi}}^N(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}' \in \text{supp}(\hat{\rho}_T^N)} \|\boldsymbol{\phi}(\boldsymbol{\theta}) - \boldsymbol{\phi}(\boldsymbol{\theta}')\|.$$

We denote $\bar{\theta}_{i_\theta} \in \text{supp}(\hat{\rho}_T^N)$ such that $\bar{\phi}^N(\theta) = \phi(\bar{\theta}_{i_\theta})$. It implies that

$$\begin{aligned}
 \sup_{\theta \in \text{supp}(\rho_T^\infty)} \|\phi(\theta) - \bar{\phi}^N(\theta)\| &\leq \sup_{\theta \in \text{supp}(\rho_T^\infty)} \|\phi(\theta) - \phi(\theta_{i_\theta})\| \\
 &= \sup_{\theta \in \text{supp}(\rho_T^\infty)} \|\phi(\theta) - \bar{\phi}^N(\theta) + \bar{\phi}^N(\theta) - \phi(\theta_{i_\theta})\| \\
 &= \sup_{\theta \in \text{supp}(\rho_T^\infty)} \|\phi(\theta) - \phi(\bar{\theta}_{i_\theta}) + \phi(\bar{\theta}_{i_\theta}) - \phi(\theta_{i_\theta})\| \\
 &\leq \sup_{\theta \in \text{supp}(\rho_T^\infty)} \|\phi(\theta) - \phi(\bar{\theta}_{i_\theta})\| + \sup_{\theta \in \text{supp}(\rho_T^\infty)} \|\phi(\theta) - \phi(\bar{\theta}_{i_\theta})\| \\
 &\leq \epsilon/2 + \max_{i \in [N]} \|\bar{\theta}_i(T) - \theta_i(T)\| \|\phi\|_{\text{Lip}} \\
 &\leq \epsilon.
 \end{aligned}$$