

---

# A Theory-Driven Self-Labeling Refinement Method for Contrastive Representation Learning (Supplementary File)

---

Pan Zhou\* Caiming Xiong\* Xiao-Tong Yuan† Steven Hoi\*

\* Salesforce Research

† Nanjing University of Information Science & Technology  
panzhou3@gmail.com {cxiong, shoi}@salesforce.com xtyuan@nuist.edu.cn

## Abstract

This supplementary document contains more additional experimental details and the technical proofs of convergence results of the NeurIPS’21 submission entitled “A Theory-Driven Self-Labeling Refinement Method for Contrastive Representation Learning”. It is structured as follows. In Appendix A, we provide more experimental details, including training algorithm, network architecture, optimizer details, loss construction and training cost of SANE. Appendix B presents the proof and details of the main results, namely, Theorem 1, in Section 2, which analyzes the generalization performance of MoCo.

Next, Appendix C introduces the *proof roadmap* and details of the main results, i.e. Theorem 2, in Section 3.1. Since the proof framework is relatively complex, we first introduce some necessary preliminaries, including notations, conceptions and assumptions that are verified in subsequent analysis in Appendix C.2.4. Then we provide the proofs of Theorem 2 in Appendix C.2. Specifically, we first introduce the *proof roadmap of Theorem 2* in Appendix C.2.1. Then we present several auxiliary theories in Appendix C.2.2. Next, we prove our Theorem 2 in Appendix C.2.3. Finally, we present all proof details of auxiliary theories in Appendix C.2.4.

## A More Experimental Details

Due to space limitation, we defer more experimental details to this appendix. Here we first introduce the training algorithm of SANE, and then present more setting details of optimizers, architectures, loss construction for CIFAR10 and ImageNet.

### A.1 Algorithm Framework of SANE

In this subsection, we introduce the training algorithm of SANE in details, which is summarized in Algorithm 1. Same as MoCo [1] and CLSA [2], we alternatively update the online network  $f_w$  and target network  $g_\xi$  via SGD optimizer. Our codes are implemented based on MoCo and CLSA. The code of MoCo and CLSA satisfies “Creative Commons Attribution-NonCommercial 4.0 International Public License”.

---

**Algorithm 1** Algorithm Framework for SANE

---

**Input:** online network  $f_w$ , target network  $g_\xi$ , dictionary  $B$ , temperature parameter  $\tau$ , momentum-update parameter  $\iota$ , sharpness parameter  $\tau'$ , prior confidence  $\mu$ , regularization weight  $\lambda$ , parameter  $\kappa$  for  $\text{Beta}(\kappa, \kappa)$ , weak augmentation  $T_1$ , and weak or strong augmentation  $T_2$

**Initialization:** initialize online network  $f_w$ , target network  $g_\xi$ , dictionary  $B$  as MoCo.

**for**  $i = 1 \dots T$  **do**

1. sample a minibatch of vanilla samples  $\{c_i\}_{i=1}^s$
2. use  $T_1$  to augment  $\{c_i\}_{i=1}^s$  to obtain weak augmentations  $\{(x_i, \tilde{x}_i)\}_{i=1}^s$ , i.e.  $x_i = T_1(c_i)$  and  $\tilde{x}_i = T_1(c_i)$ .
3. compute feature  $\{f(x_i)\}_{i=1}^s$  and  $B' = \{g(\tilde{x}_i)\}_{i=1}^s$
4. compute the contrastive loss  $\mathcal{L}_c(w, \{(x_i, y_i)\})$  in Eqn. (9)
5. use  $\tilde{x}_i$  to compute the estimated labels  $\tilde{y}_i^t$  of query  $x_i$  by self-labeling refinery (5) ( $\forall i = 1, \dots, s$ )
6. if using strong augmentation for momentum mixup, use  $T_2$  to augment  $\{c_i\}_{i=1}^s$  for obtaining strong augmentations  $\{\tilde{x}_i\}_{i=1}^s$  to replace the previous  $\{\tilde{x}_i\}_{i=1}^s$  in  $\{(x_i, \tilde{x}_i)\}_{i=1}^s$
7. use momentum mixup (8) and samples  $\{(x_i, \tilde{x}_i, \tilde{y}_i^t)\}_{i=1}^s$  to obtain new virtual queries and labels  $\{(x'_i, y'_i)\}_{i=1}^s$
8. use  $\{(x'_i, y'_i)\}_{i=1}^s$  to compute the momentum mixup contrastive loss  $\mathcal{L}_c(w, \{(x'_i, y'_i)\})$  in Eqn. (9)
9. update online network  $f_w$  by minimizing  $(1-\lambda)\mathcal{L}_c(w, \{(x_i, y_i)\}) + \lambda\mathcal{L}_c(w, \{(x'_i, y'_i)\})$
10. update target network  $g_\xi$  by exponential moving average
11. update the dictionary  $B$  via minibatch feature  $B'$  in a first-in first-out order.

**end for**

**Output:**

---

## A.2 Algorithm Parameter Settings

**Experimental Settings for Linear Evaluation on CIFAR10 and ImageNet.** For CIFAR10 and ImageNet, we follow [1, 3] and use ResNet50 [4] as a backbone. Then we first pretrain SANE on the corresponding training data, and then train a linear classifier on top of 2048-dimensional frozen features provided by ResNet50. For pretraining on both datasets, we use SGD with an initial learning rate 0.03 (annealed down to zero via cosine decay [5]), a momentum of 0.9, and a weight decay of  $10^{-4}$ . Such optimizer parameters are the same with MoCo and CLSA.

Next, we pretrain 2,000 epochs on CIFAR10 with minibatch size 256 and dictionary size 4,096. For pretraining on Imagenet, the dictionary size is always 65,536; the batch size is often 256 on a cluster of 8 GPUs and is linearly scaled together with learning rate on multiple clusters. For linear classifier training, we use ADAM [6] with a learning rate of 0.01 and without weight decay to train 200 epochs on CIFAR10, and adopt SGD with an initial learning 10 (cosine decayed to zero) and a momentum of 0.9 to train 100 epochs on ImageNet. We use standard data augmentations in [1] for pretraining unless otherwise stated. Specifically, for pretraining on CIFAR10 and ImageNet, we follow MoCo and use RandomResizedCrop, ColorJitter, RandomGrayscale, GaussianBlur, RandomHorizontalFlip, and Normalization. For CIFAR10, please find its pretraining augmentation in the example<sup>1</sup>. Except the above random augmentation, we also use the proposed momentum mixup to generate the virtual instances for constructing the momentum mixup loss.

For CIFAR10, to fairly compare with [7], we crop each image into two views to construct the loss (9). Specifically, for a minibatch of vanilla samples  $\{c_i\}_{i=1}^s$ , we use weak augmentation  $T_1$  to augment  $\{c_i\}_{i=1}^s$  to obtain weak augmentations  $\{(x_i, \tilde{x}_i)\}_{i=1}^s$ , i.e.  $x_i = T_1(c_i)$  and  $\tilde{x}_i = T_1(c_i)$ . Then same as MoCo, we can compute the contrastive loss by using  $\{(x_i, \tilde{x}_i)\}_{i=1}^s$ . Meanwhile, we use  $\tilde{x}_i$  to compute the soft label  $\tilde{y}_i^t$  of  $\tilde{x}_i$  via (5). Next, we use momentum mixup (8) and samples  $\{(x_i, \tilde{x}_i, \tilde{y}_i^t)\}_{i=1}^s$  to obtain new virtual queries and labels  $\{(x'_i, y'_i)\}_{i=1}^s$ , and then use  $\{(x'_i, y'_i)\}_{i=1}^s$  to compute the momentum mixup contrastive loss  $\mathcal{L}_c(w, \{(x'_i, y'_i)\})$  in Eqn. (9). For strong augmentation, after we compute the vanilla contrastive loss in MoCo, and then use strong augmentation to augment  $\{c_i\}_{i=1}^s$  to replace  $\tilde{x}_i$  in  $\{(x_i, \tilde{x}_i, \tilde{y}_i^t)\}_{i=1}^s$ . Then we can generate virtual

---

<sup>1</sup>[https://colab.research.google.com/github/facebookresearch/moco/blob/colab-notebook/colab/moco\\_cifar10\\_demo.ipynb](https://colab.research.google.com/github/facebookresearch/moco/blob/colab-notebook/colab/moco_cifar10_demo.ipynb)

query instances and their labels ( $\{(\mathbf{x}'_i, \mathbf{y}'_i)\}_{i=1}^s$ ) by using  $\{(\mathbf{x}_i, \tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i^t)\}_{i=1}^s$ . The training cost on CIFAR10 for 2,000 epochs is about 11 days on single V100 GPU.

For ImageNet, we follow CLSA for fair comparison. For SANE-Single, we use the same way to construct the contrastive loss, and then use augmentation  $T_1$  to augment  $\{c_i\}_{i=1}^s$  to replace  $\tilde{\mathbf{x}}_i$  in  $\{(\mathbf{x}_i, \tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i^t)\}_{i=1}^s$  to construct the momentum mixup loss. Indeed, we also can do not replace  $\tilde{\mathbf{x}}_i$  in  $\{(\mathbf{x}_i, \tilde{\mathbf{x}}_i, \tilde{\mathbf{y}}_i^t)\}_{i=1}^s$  for momentum mixup loss, which actually did not affect the performance. We do it, since SANE-Multi crops each image into five different crops for constructing momentum mixup loss, and thus SANE-Single and SANE-Multi will be more consistent, i.e. SANE-Multi uses 5 crops while SANE-Single uses one crop. For strong augmentation, we replace the augmentation  $T_1$  in momentum mixup with strong augmentation, which is the same on CIFAR10. As mentioned above, to construct the momentum mixup loss, SANE-Multi crops each image into five sizes  $224 \times 224$ ,  $192 \times 192$ ,  $160 \times 160$ ,  $128 \times 128$ , and  $96 \times 96$  and averages their momentum mixup losses. For the vanilla contrastive loss, SANE-Multi uses the same way in SANE-Single to compute. In this way, SANE-Single and SANE-Multi respectively have the same settings with CLSA-Single and CLSA-Multi. Thus, ELSE has almost the same training cost with CLSA, i.e. about 75 (188) hours with 8 GPUs, 200 epochs, batch size of 256 for SANE-Single (-Multi). It should be mentioned that for vanilla contrastive loss in both CLSA-Single and CLSA-Multi, we always use weak augmentations.

**Transfer Evaluation Settings.** We evaluate the pretrained model on ImageNet on VOC [8] and COCO [9]. For VOC, similar to linear evaluation, we train a linear classifier upon ResNet50 100 epochs by SGD with a learning rate 0.05, a momentum 0.9, batch size 256, and without weight and learning rate decay. For COCO, we adopt the same protocol in [1] to fine-tune the pretrained ResNet50 based on detectron2 [10] for fairness. We evaluate the transfer ability of the cells selected on CIFAR10 by testing them on ImageNet. Following DARTS, we use momentum SGD with an initial learning 0.025 (cosine decayed to zero), a momentum of 0.9, a weight decay of  $3 \times 10^{-4}$ , and gradient norm clipping parameter 5.0.

## B Proofs of The Results in Section 2

**Lemma 1.** [11] *Suppose the loss  $\ell$  is bounded by the range  $[a, b]$ , namely  $\ell(f(\mathbf{x}; \mathbf{w}), \mathbf{y}) \in [a, b]$ . Then let  $\mathcal{F}$  be a finite class of hypotheses  $\ell(f(\mathbf{x}; \mathbf{w}), \mathbf{y}) : \mathcal{X} \rightarrow \mathbb{R}$ . Let*

$$\mathcal{Q}_e(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i; \mathbf{w}), \mathbf{y}_i), \quad \mathcal{Q}(f) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \in \mathcal{S}} [\ell(f(\mathbf{x}; \mathbf{w}), \mathbf{y})]$$

*respectively denote the empirical and population risk, where  $\mathcal{S}$  denote the unknown data distribution and the sampled dataset  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n \sim \mathcal{S}$  is of size  $n$ . Then for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$  we have*

$$\mathcal{Q}(f) \leq \mathcal{Q}_e(f) + \sqrt{\frac{2(b-a)^2 V_{\mathcal{D}} \ln(2|\mathcal{F}|/\delta)}{n}} + \frac{7(b-a)^2 \ln(2|\mathcal{F}|/\delta)}{3(n-1)}, \quad (7)$$

*where  $V_{\mathcal{D}}$  denotes the variance of the loss  $\ell(f(\mathbf{x}; \mathbf{w}), \mathbf{y})$  on the dataset  $\mathcal{D}$ , and  $|\mathcal{F}|$  denotes the covering number of  $\mathcal{F}$  in the uniform norm  $\|\cdot\|_{\infty}$ .*

**Lemma 2.** [12] *For any polynomials  $f(x) = \sum_{i=0}^p a_i x^i$ ,  $x \in [0, 1]$  and  $\sum_{i=1}^p |a_i| < 1$ , there exists a multilayer neural network  $\hat{f}(x)$  with  $\mathcal{O}(p + \log \frac{p}{\epsilon})$  layers,  $\mathcal{O}(\log \frac{p}{\epsilon})$  binary step units and  $\mathcal{O}(p \log \frac{p}{\epsilon})$  rectifier linear units such that  $|f(x) - \hat{f}(x)| \leq \epsilon$ ,  $\forall x \in [0, 1]$ .*

*Assume that function  $f$  is continuous on  $[0, 1]$  and  $\lceil \log \frac{2}{\epsilon} \rceil + 1$  times differential in  $(0, 1)$ . Let  $f^{(n)}$  denote the derivative of  $f$  of  $n$ -th order and  $\|f\| = \max_{x \in [0, 1]} |f(x)|$ . If  $\|f^{(n)}\| \leq n!$  holds for all  $n \in [\lceil \log \frac{2}{\epsilon} \rceil + 1, \infty)$ , then there exists a deep network  $f$  with  $\mathcal{O}(\log \frac{1}{\epsilon})$  layers,  $\mathcal{O}(\log \frac{1}{\epsilon})$  binary step units and  $\mathcal{O}(\log^2 \frac{1}{\epsilon})$  rectifier linear units such that  $|f(x) - \hat{f}(x)| \leq \epsilon$ ,  $\forall x \in [0, 1]$ .*

For expression power analysis of deep network, more stronger results can be found in [13, 14, 15, 16] and all show that any function can be approximately can be approximated by a deep network to arbitrary accuracy.

## B.1 Proof of Theorem 1

*Proof.* Here we use two steps to prove our results in Theorem 1.

$$\tilde{\mathcal{Q}}(f_w) = \frac{1}{n} \sum_{i=1}^n \ell(h(f_w(\mathbf{x}_i), \mathbf{B}_i), \mathbf{y}_i), \quad (8)$$

**Step 1. proof for first part results.** To begin with, we first define an empirical risk  $\mathcal{Q}_e(f)$ :

$$\mathcal{Q}_e(f) = \frac{1}{n} \sum_{i=1}^n \ell(h(f_w(\mathbf{x}_i), \mathbf{B}_i), \mathbf{y}_i^*),$$

where  $\mathcal{Q}_e(f)$  uses the ground truth label  $\mathbf{y}_i^*$  for training. From Lemma 1, with probability at least  $1 - \delta$ , we have

$$\mathcal{Q}(f) \leq \mathcal{Q}_e(f) + \sqrt{\frac{2(b-a)^2 V_{\mathcal{D}} \ln(2|\mathcal{F}|/\delta)}{n}} + \frac{7(b-a)^2 \ln(2|\mathcal{F}|/\delta)}{3(n-1)},$$

where  $\mathcal{Q}(f)$  is the population risk, and  $\mathcal{Q}_e(f)$  is the empirical risk. Both are trained with the ground truth  $\mathbf{y}_i^*$ . So the remaining work is to upper bound  $\mathcal{Q}_e(f)$  via  $\tilde{\mathcal{Q}}(f)$ . Towards this end, we can bound it as follows

$$\begin{aligned} \mathcal{Q}_e(f) - \tilde{\mathcal{Q}}(f) &= \frac{1}{n} \sum_{i=1}^n (\ell(h(f_w(\mathbf{x}_i), \mathbf{B}_i), \mathbf{y}_i^*) - \ell(h(f_w(\mathbf{x}_i), \mathbf{B}_i), \mathbf{y}_i)) \\ &\stackrel{\textcircled{1}}{\leq} \frac{1}{n} \sum_{i=1}^n \|\nabla_{\mathbf{y}} \ell(h(f_w(\mathbf{x}_i), \mathbf{B}_i), \mathbf{y})\| \cdot \|\mathbf{y}_i^* - \mathbf{y}_i\|_2 \\ &\stackrel{\textcircled{2}}{\leq} L_y \mathbb{E}_i \|\mathbf{y}_i^* - \mathbf{y}_i\|_2 \\ &\stackrel{\textcircled{2}}{\leq} L_y \mathbb{E}_{\mathcal{D} \sim \mathcal{S}} [\|\mathbf{y}^* - \mathbf{y}\|_2], \end{aligned}$$

where  $\textcircled{1}$  holds by using  $\mathbf{y} = \mathbf{y}_i + \theta(\mathbf{y}_i^* - \mathbf{y}_i)$  for certain  $\theta \in (0, 1)$ ;  $\textcircled{2}$  holds since we use the  $L_y$ -Lipschitz property of  $\ell(h(f_w(\mathbf{x}_i), \mathbf{B}_i), \mathbf{y}_i)$ . Then combining these results together, we can obtain the desired results:

$$|\mathcal{Q}(f) - \tilde{\mathcal{Q}}(f)| \leq L_y \mathbb{E}_{\mathcal{D} \sim \mathcal{S}} \|\mathbf{y}_i^* - \mathbf{y}_i\|_2 + \sqrt{\frac{2(b-a)^2 V_{\mathcal{D}} \ln(2|\mathcal{F}|/\delta)}{n}} + \frac{7(b-a)^2 \ln(2|\mathcal{F}|/\delta)}{3(n-1)}.$$

**Step 2. proof for second part results.** Here we can construct a simple two-classification problem for clarity. Suppose we have two classes: class one with training data  $\mathcal{D}_1 = \{(\mathbf{x}_1, \mathbf{x}_1, \mathbf{y}_1^*)\}_{i=1}^{n/2}$  and class two with training data  $\mathcal{D}_2 = \{(\mathbf{x}_2, \mathbf{x}_2, \mathbf{y}_2^*)\}_{i=1}^{n/2}$ , where  $\mathbf{y}_1^*$  denotes the ground truth label of  $\mathbf{x}_1$  on the set  $\mathbf{B}_1 = \{\mathbf{x}_1 \cup \mathbf{B}\}$ , and  $\mathbf{y}_2^*$  denotes the ground truth label of  $\mathbf{x}_2$  on the set  $\mathbf{B}_2 = \{\mathbf{x}_2 \cup \mathbf{B}\}$ . Both training datasets  $\mathcal{D}_1$  and  $\mathcal{D}_2$  have  $\frac{n}{2}$  samples. Here we assume there is no data augmentation which means  $\mathbf{x}_i = \tilde{\mathbf{x}}_i$  in the manuscript. In  $\mathcal{D}_1$ , its samples are the same, namely  $(\mathbf{x}_1, \mathbf{x}_1, \mathbf{y}_1^*)$ . Similarly,  $\mathcal{D}_2$  also has the same samples, namely  $(\mathbf{x}_2, \mathbf{x}_2, \mathbf{y}_2^*)$ . Then the predicted class probability  $\mathbf{y}_{ij}$  of sample  $\mathbf{x}_i$  on class  $j$  is as follows:

$$\mathbf{y}_{i0} = \frac{e^{\delta(\mathbf{x}_i, \mathbf{x}_i)/t}}{e^{\delta(\mathbf{x}_i, \mathbf{x}_i)/\tau} + \sum_{j=1}^k e^{\delta(\mathbf{x}_i, \mathbf{b}_j)/\tau}}, \quad \mathbf{y}_{ij} = \frac{e^{\delta(\mathbf{x}_i, \mathbf{b}_j)/\tau}}{e^{\delta(\mathbf{x}_i, \mathbf{x}_i)/\tau} + \sum_{j=1}^k e^{\delta(\mathbf{x}_i, \mathbf{b}_j)/\tau}} \quad (j = 1, \dots, k), \quad (9)$$

where  $\delta(\mathbf{x}_i, \tilde{\mathbf{x}}_i) = -\frac{\langle f(\mathbf{x}_i), g(\tilde{\mathbf{x}}_i) \rangle}{\|f(\mathbf{x}_i)\|_2 \cdot \|g(\tilde{\mathbf{x}}_i)\|_2}$ ,  $\tau$  denotes a temperature. For simplicity, we let dictionary  $\mathbf{B} = \{\mathbf{x}_1, \mathbf{x}_2\}$ . In this way, we have for both ground truth label  $\mathbf{y}_1^*$  and  $\mathbf{y}_2^*$  that satisfy  $\mathbf{y}_{10}^* = \mathbf{y}_{11}^*$ ,  $\mathbf{y}_{10}^* + \mathbf{y}_{11}^* + \mathbf{y}_{12}^* = 1$ ,  $\mathbf{y}_{20}^* = \mathbf{y}_{22}^*$ ,  $\mathbf{y}_{20}^* + \mathbf{y}_{21}^* + \mathbf{y}_{22}^* = 1$ . For this setting, here we assume the training labels are denoted by  $\mathbf{y}_1$  and  $\mathbf{y}_2$ . Moreover, they satisfy  $\mathbf{y}_{10} = \mathbf{y}_{11} > 0$ ,  $\mathbf{y}_{10} + \mathbf{y}_{11} + \mathbf{y}_{12} = 1$ ,  $\mathbf{y}_{20} = \mathbf{y}_{22} > 0$ ,  $\mathbf{y}_{20} + \mathbf{y}_{21} + \mathbf{y}_{22} = 1$ . The reason that we do not use one-hot labels. This is because for dictionary  $\mathbf{B} = \{\mathbf{x}_1, \mathbf{x}_2\}$ , given a sample  $\mathbf{x}_i$  ( $i = 1, 2$ ),  $\mathbf{x}_i$  needs to predict the labels on the set  $\{\mathbf{x}_i \cup \mathbf{B}\} = \{\mathbf{x}_i, \mathbf{x}_1, \mathbf{x}_2\}$ , where the labels are not one-hot obviously and satisfy  $\mathbf{y}_{i1} = \mathbf{y}_{ii} > 0$ . In the following, we will train the model on the training data  $\tilde{\mathcal{D}} = \tilde{\mathcal{D}}_1 \cup \tilde{\mathcal{D}}_2$  where  $\tilde{\mathcal{D}}_1 = \{(\mathbf{x}_1, \mathbf{x}_1, \mathbf{y}_1)\}$  and  $\tilde{\mathcal{D}}_2 = \{(\mathbf{x}_2, \mathbf{x}_2, \mathbf{y}_2)\}$ . We use  $\tilde{\mathbf{y}}_i$  to denote the model predicted label of  $\mathbf{x}_i$ .

Then for the test samples, we assume that half of samples are  $(\mathbf{x}_1, \mathbf{x}_1, \mathbf{y}_1^*)$  and remaining samples are  $(\mathbf{x}_2, \mathbf{x}_2, \mathbf{y}_2^*)$ . Then for any network  $f$ , we always have

$$\mathcal{Q}(f) - \mathcal{Q}_e(f) = \frac{1}{n} \sum_{i=1}^n (\mathbb{E}[\ell(h(f_{\mathbf{w}}(\mathbf{x}_i), \mathbf{B}_i), \mathbf{y}_i^*)] - \ell(h(f_{\mathbf{w}}(\mathbf{x}_i), \mathbf{B}_i), \mathbf{y}_i^*)) = 0.$$

Then we attempt to lower bound  $\mathcal{Q}_e(f) - \tilde{\mathcal{Q}}(f)$ . Our training dataset is  $\tilde{\mathcal{D}} = \tilde{\mathcal{D}}_1 \cup \tilde{\mathcal{D}}_2$  where  $\tilde{\mathcal{D}}_1 = \{(\mathbf{x}_1, \mathbf{x}_1, \mathbf{y}_1)\}$  and  $\tilde{\mathcal{D}}_2 = \{(\mathbf{x}_2, \mathbf{x}_2, \mathbf{y}_2)\}$ . Then we discuss whether the network  $f$  can perfectly fit the labels (9) of data  $\tilde{\mathcal{D}}$ . For both cases, our results can hold.

**Perfectly fitting.** Network  $f$  has the capacity to perfectly fit the label  $\tilde{\mathbf{y}}_1$  in  $\tilde{\mathcal{D}}_1$  and the label  $\tilde{\mathbf{y}}_2$  in  $\tilde{\mathcal{D}}_2$  when  $\mathbf{x}_1$  are different  $\mathbf{x}_2$ . In this case, we have

$$\begin{aligned} & \mathcal{Q}_e(f) - \tilde{\mathcal{Q}}(f) \\ &= \frac{1}{n} \sum_{i=1}^n (\ell(h(f_{\mathbf{w}}(\mathbf{x}_i), \mathbf{B}_i), \mathbf{y}_i^*) - \ell(h(f_{\mathbf{w}}(\mathbf{x}_i), \mathbf{B}_i), \mathbf{y}_i)) \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^k (\mathbf{y}_{i,s}^* \log(h(f_{\mathbf{w}}(\mathbf{x}_i), \mathbf{B}_i)) - \mathbf{y}_{i,s} \log(h(f_{\mathbf{w}}(\mathbf{x}_i), \mathbf{B}_i))) \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^k (\mathbf{y}_{i,s}^* - \mathbf{y}_{i,s}) \log(\tilde{\mathbf{y}}_{i,s}) \\ &\stackrel{\textcircled{1}}{=} \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^k (\mathbf{y}_{i,s}^* - \mathbf{y}_{i,s}) \log(\mathbf{y}_{i,s}) \\ &= \frac{1}{6} [(\mathbf{y}_{10}^* - \mathbf{y}_{10}) \log(\mathbf{y}_{10}) + (\mathbf{y}_{11}^* - \mathbf{y}_{11}) \log(\mathbf{y}_{11}) + (\mathbf{y}_{12}^* - \mathbf{y}_{12}) \log(\mathbf{y}_{12}) \\ &\quad + (\mathbf{y}_{20}^* - \mathbf{y}_{20}) \log(\mathbf{y}_{20}) + (\mathbf{y}_{21}^* - \mathbf{y}_{21}) \log(\mathbf{y}_{21}) + (\mathbf{y}_{22}^* - \mathbf{y}_{22}) \log(\mathbf{y}_{22})] \\ &= \frac{1}{6} [2(\mathbf{y}_{10}^* - \mathbf{y}_{10}) \log(\mathbf{y}_{10}) + (\mathbf{y}_{12}^* - \mathbf{y}_{12}) \log(\mathbf{y}_{12}) + 2(\mathbf{y}_{20}^* - \mathbf{y}_{20}) \log(\mathbf{y}_{20}) + (\mathbf{y}_{21}^* - \mathbf{y}_{21}) \log(\mathbf{y}_{21})] \\ &\stackrel{\textcircled{2}}{=} \frac{1}{3} \left[ (\mathbf{y}_{10}^* - \mathbf{y}_{10}) \log \frac{\mathbf{y}_{10}}{1 - 2\mathbf{y}_{10}} + (\mathbf{y}_{20}^* - \mathbf{y}_{20}) \log \frac{\mathbf{y}_{20}}{1 - 2\mathbf{y}_{20}} \right], \end{aligned}$$

where  $\textcircled{1}$  holds since  $\tilde{\mathbf{y}}_{i,s} = \mathbf{y}_{i,s}$ , and  $\textcircled{2}$  uses  $\mathbf{y}_{10}^* = \mathbf{y}_{11}^*$ ,  $\mathbf{y}_{10}^* + \mathbf{y}_{11}^* + \mathbf{y}_{12}^* = 1$ ,  $\mathbf{y}_{20}^* = \mathbf{y}_{22}^*$ ,  $\mathbf{y}_{20}^* + \mathbf{y}_{21}^* + \mathbf{y}_{22}^* = 1$ ,  $\mathbf{y}_{10} = \mathbf{y}_{11}$ ,  $\mathbf{y}_{10} + \mathbf{y}_{11} + \mathbf{y}_{12} = 1$ ,  $\mathbf{y}_{20} = \mathbf{y}_{22}$ ,  $\mathbf{y}_{20} + \mathbf{y}_{21} + \mathbf{y}_{22} = 1$ . Then we can choose proper values such that

$$\mathbf{y}_{10}^* = \mathbf{y}_{11}^* > \mathbf{y}_{10} = \mathbf{y}_{11} > \frac{1}{3}, \mathbf{y}_{20}^* = \mathbf{y}_{22}^* > \mathbf{y}_{20} = \mathbf{y}_{22} > \frac{1}{3}.$$

For example, we can let  $\mathbf{y}_1 = (0.4, 0.4, 0.2)$ ,  $\mathbf{y}_1^* = (0.45, 0.45, 0.1)$ ,  $\mathbf{y}_2 = (0.4, 0.2, 0.4)$ ,  $\mathbf{y}_2^* = (0.45, 0.1, 0.45)$ . In this way, we have  $(\mathbf{y}_{10}^* - \mathbf{y}_{10}) \log \frac{\mathbf{y}_{10}}{1 - 2\mathbf{y}_{10}} \geq c_1(\mathbf{y}_{10}^* - \mathbf{y}_{10}) > 0$  and  $(\mathbf{y}_{20}^* - \mathbf{y}_{20}) \log \frac{\mathbf{y}_{20}}{1 - 2\mathbf{y}_{20}} \geq c_2(\mathbf{y}_{20}^* - \mathbf{y}_{20}) > 0$ . So this means that there exists a constant  $C$  such that

$$\mathcal{Q}_e(f) - \tilde{\mathcal{Q}}(f) \geq C \cdot \mathbb{E}_i [\|\mathbf{y}_i^* - \mathbf{y}_i\|_2] = C \cdot \mathbb{E}_{\mathcal{D} \sim \mathcal{S}} [\|\mathbf{y}^* - \mathbf{y}\|_2].$$

So combining the above results gives the following desired result:

$$\mathcal{Q}(f) - \tilde{\mathcal{Q}}(f) \geq C \cdot \mathbb{E}_{\mathcal{D} \sim \mathcal{S}} [\|\mathbf{y}^* - \mathbf{y}\|_2].$$

**Non-perfectly fitting.** From Lemma 2 (other more results in [13, 14, 15, 16]), one can approximate any function by a deep network to arbitrary accuracy. Specifically, for the polynomial function in Eqn. (9), there exists a multilayer neural network  $\hat{f}(x)$  with proper width and depth such that  $\|\mathbf{y}_1 - \tilde{\mathbf{y}}_1\|_1 \leq \epsilon$  and  $\|\mathbf{y}_2 - \tilde{\mathbf{y}}_2\|_1 \leq \epsilon$ , where  $\tilde{\mathbf{y}}_1$  and  $\tilde{\mathbf{y}}_2$  are the predicted labels of samples  $\mathbf{x}_1$  and  $\mathbf{x}_2$  by using (9). The labels  $\mathbf{y}_1$  and  $\mathbf{y}_2$  are associated with our training dataset  $\tilde{\mathcal{D}} = \tilde{\mathcal{D}}_1 \cup \tilde{\mathcal{D}}_2$  where

$\tilde{\mathcal{D}}_1 = \{(\mathbf{x}_1, \mathbf{x}_1, \mathbf{y}_1)\}$  and  $\tilde{\mathcal{D}}_2 = \{(\mathbf{x}_2, \mathbf{x}_2, \mathbf{y}_2)\}$ . In this case, we have

$$\begin{aligned}
& \mathcal{Q}_e(f) - \tilde{\mathcal{Q}}(f) \\
&= \frac{1}{n} \sum_{i=1}^n (\ell(h(f_{\mathbf{w}}(\mathbf{x}_i), \mathbf{B}_i), \mathbf{y}_i^*) - \ell(h(f_{\mathbf{w}}(\mathbf{x}_i), \mathbf{B}_i), \mathbf{y}_i)) \\
&= \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^k (\mathbf{y}_{i,s}^* \log(h(f_{\mathbf{w}}(\mathbf{x}_i), \mathbf{B}_i)) - \mathbf{y}_{i,s} \log(h(f_{\mathbf{w}}(\mathbf{x}_i), \mathbf{B}_i))) \\
&= \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^k (\mathbf{y}_{i,s}^* - \mathbf{y}_{i,s}) \log(\tilde{\mathbf{y}}_{i,s}) \\
&= \frac{1}{6} [(\mathbf{y}_{10}^* - \mathbf{y}_{10}) \log(\tilde{\mathbf{y}}_{10}) + (\mathbf{y}_{11}^* - \mathbf{y}_{11}) \log(\tilde{\mathbf{y}}_{11}) + (\mathbf{y}_{12}^* - \mathbf{y}_{12}) \log(\tilde{\mathbf{y}}_{12}) \\
&\quad + (\mathbf{y}_{20}^* - \mathbf{y}_{20}) \log(\tilde{\mathbf{y}}_{20}) + (\mathbf{y}_{21}^* - \mathbf{y}_{21}) \log(\tilde{\mathbf{y}}_{21}) + (\mathbf{y}_{22}^* - \mathbf{y}_{22}) \log(\tilde{\mathbf{y}}_{22})] \\
&= \frac{1}{6} [2(\mathbf{y}_{10}^* - \mathbf{y}_{10}) \log(\tilde{\mathbf{y}}_{10}) + (\mathbf{y}_{12}^* - \mathbf{y}_{12}) \log(\tilde{\mathbf{y}}_{12}) + 2(\mathbf{y}_{20}^* - \mathbf{y}_{20}) \log(\tilde{\mathbf{y}}_{20}) + (\mathbf{y}_{21}^* - \mathbf{y}_{21}) \log(\tilde{\mathbf{y}}_{21})] \\
&\stackrel{\textcircled{1}}{=} \frac{1}{3} \left[ (\mathbf{y}_{10}^* - \mathbf{y}_{10}) \log \frac{\tilde{\mathbf{y}}_{10}}{1 - 2\tilde{\mathbf{y}}_{10}} + (\mathbf{y}_{20}^* - \mathbf{y}_{20}) \log \frac{\tilde{\mathbf{y}}_{20}}{1 - 2\tilde{\mathbf{y}}_{20}} \right],
\end{aligned}$$

where  $\textcircled{1}$  uses  $\mathbf{y}_{10}^* = \mathbf{y}_{11}^*$ ,  $\mathbf{y}_{10}^* + \mathbf{y}_{11}^* + \mathbf{y}_{12}^* = 1$ ,  $\mathbf{y}_{20}^* = \mathbf{y}_{21}^*$ ,  $\mathbf{y}_{20}^* + \mathbf{y}_{21}^* + \mathbf{y}_{22}^* = 1$ ,  $\mathbf{y}_{10} = \mathbf{y}_{11}$ ,  $\mathbf{y}_{10} + \mathbf{y}_{11} + \mathbf{y}_{12} = 1$ ,  $\mathbf{y}_{20} = \mathbf{y}_{22}$ ,  $\mathbf{y}_{20} + \mathbf{y}_{21} + \mathbf{y}_{22} = 1$ . Then we can choose proper values such that

$$\mathbf{y}_{10}^* = \mathbf{y}_{11}^* > \mathbf{y}_{10} = \mathbf{y}_{11} > \frac{1}{3} + \epsilon, \mathbf{y}_{20}^* = \mathbf{y}_{22}^* > \mathbf{y}_{20} = \mathbf{y}_{22} > \frac{1}{3} + \epsilon.$$

For example, we can let  $\mathbf{y}_1 = (0.4, 0.4, 0.2)$ ,  $\mathbf{y}_1^* = (0.45, 0.45, 0.1)$ ,  $\mathbf{y}_2 = (0.4, 0.2, 0.4)$ ,  $\mathbf{y}_2^* = (0.45, 0.1, 0.45)$ , and  $\epsilon = 0.0001$ . In this way, we have  $(\mathbf{y}_{10}^* - \mathbf{y}_{10}) \log \frac{\tilde{\mathbf{y}}_{10}}{1 - 2\tilde{\mathbf{y}}_{10}} \geq c_1(\mathbf{y}_{10}^* - \mathbf{y}_{10}) > 0$  and  $(\mathbf{y}_{20}^* - \mathbf{y}_{20}) \log \frac{\tilde{\mathbf{y}}_{20}}{1 - 2\tilde{\mathbf{y}}_{20}} \geq c_2(\mathbf{y}_{20}^* - \mathbf{y}_{20}) > 0$ . So this means that there exists a constant  $C$  such that

$$\mathcal{Q}_e(f) - \tilde{\mathcal{Q}}(f) \geq C \cdot \mathbb{E}_i [\|\mathbf{y}_i^* - \mathbf{y}_i\|_2] = C \cdot \mathbb{E}_{\mathcal{D} \sim \mathcal{S}} [\|\mathbf{y}^* - \mathbf{y}\|_2].$$

So combining the above results gives the following desired result:

$$\mathcal{Q}(f) - \tilde{\mathcal{Q}}(f) \geq C \cdot \mathbb{E}_{\mathcal{D} \sim \mathcal{S}} [\|\mathbf{y}^* - \mathbf{y}\|_2].$$

The proof is completed.  $\square$

## C Proof of Results in Section 3.1

In this section, we first introduce some necessary preliminaries, including notations, conceptions and assumptions that are verified in subsequent analysis in Appendix C.2.4. Then we provide the proofs of Theorem 2 in Appendix C.2. Specifically, we first introduce the proof roadmap in Appendix C.2.1. Then we present several auxiliary theories in Appendix C.2.2. Next, we prove our Theorem 2 in Appendix C.2.3. Finally, we present all proof details of auxiliary theories in Appendix C.2.2.

### C.1 Preliminaries

#### C.1.1 General Model Formulation

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

$$\mathcal{L}_t(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^n (\tilde{\mathbf{y}}_i^t - f(\mathbf{w}, \mathbf{x}_i))^2.$$

where  $\bar{\mathbf{y}}_i^t = (1 - \alpha_t)\mathbf{y}_i + \alpha_t \mathbf{p}^t = (1 - \alpha_t)\mathbf{y}_i + \alpha_t f(\mathbf{w}_t, \mathbf{x}_i)$  denotes the estimated label of sample  $\mathbf{x}_i$ . In Assumption 2 we assume  $\beta_t = 0$  and  $\tau' = 1$  for simplicity, since performing nonlinear mapping on network output greatly increases analysis difficulty. But we will show that even though  $\beta_t = 0$  and  $\tau' = 1$ , our refinery (5) is still sufficient to refine labels. It can also be written in the more compact form

$$\mathcal{L}_t(\mathbf{w}) = \frac{1}{2} \|f(\mathbf{w}) - \bar{\mathbf{y}}^t\|_{\ell_2}^2 \quad \text{with} \quad f(\mathbf{w}) := \begin{bmatrix} f(\mathbf{w}, \mathbf{x}_1) \\ f(\mathbf{w}, \mathbf{x}_2) \\ \vdots \\ f(\mathbf{w}, \mathbf{x}_n) \end{bmatrix}. \quad (10)$$

To solve this problem we run gradient descent iterations with a constant learning rate  $\eta$  starting from an initial point  $\mathbf{w}_0$ . These iterations take the form

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla \mathcal{L}_t(\mathbf{w}_t) \quad \text{with} \quad \nabla \mathcal{L}(\mathbf{w}) = \mathcal{J}^T(\mathbf{w}) (f(\mathbf{w}) - \bar{\mathbf{y}}^t). \quad (11)$$

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

$$\mathcal{J}(\mathbf{w}) = \left[ \frac{\partial f(\mathbf{w}, \mathbf{x}_1)}{\partial \mathbf{w}} \quad \dots \quad \frac{\partial f(\mathbf{w}, \mathbf{x}_n)}{\partial \mathbf{w}} \right]^T. \quad (12)$$

Define the  $n$ -dimensional residual vector and corrupted residual vector  $\mathbf{e}$  where

$$\mathbf{r}_t = \mathbf{r}_t(\mathbf{w}) = [f(\mathbf{x}_1, \mathbf{w}_t) - \bar{\mathbf{y}}_1^t \quad \dots \quad f(\mathbf{x}_n, \mathbf{w}_t) - \bar{\mathbf{y}}_n^t]^T \quad \text{and} \quad \mathbf{e}_t = \bar{\mathbf{y}}^t - \mathbf{y}^*.$$

A key idea in our approach is that we argue that (1) in the absence of any corruption  $\mathbf{r}(\mathbf{w})$  approximately lies on the subspace  $\mathcal{S}_+$  and (2) if the labels are corrupted by a vector  $\mathbf{e}$ , then  $\mathbf{e}$  approximately lies on the complement space.

Throughout,  $\sigma_{\min}(\cdot)$  denotes the smallest singular value of a given matrix. We first introduce helpful definitions that will be used in our proofs. Given a matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and a subspace  $\mathcal{S} \subset \mathbb{R}^n$ , we define the minimum singular value of the matrix over this subspace by  $\sigma_{\min}(\mathbf{X}, \mathcal{S})$  which is defined as

$$\sigma_{\min}(\mathbf{X}, \mathcal{S}) = \sup_{\|\mathbf{v}\|_2=1, \mathbf{U}\mathbf{U}^T = \mathcal{P}_{\mathcal{S}}} \|\mathbf{v}^T \mathbf{U}^T \mathbf{X}\|_2.$$

Here,  $\mathcal{P}_{\mathcal{S}} \in \mathbb{R}^{n \times n}$  is the projection operator to the subspace. Hence, this definition essentially projects the matrix on  $\mathcal{S}$  and then takes the minimum singular value over that projected subspace.

Since augmentations are produced by using the vanilla sample  $\mathbf{c}_i$  and the augmentation  $\mathbf{x}$  obeys  $\|\mathbf{x} - \mathbf{c}_i\|_2 \leq \epsilon_0$ . So in this sense, we often call the vanilla sample and its augmentations as cluster, and call the vanilla sample as cluster center.

### C.1.2 Definitions and Assumptions

To begin with, we define  $(\epsilon, \delta)$ -clusterable dataset. As aforementioned, we often call the vanilla sample and its augmentations as cluster, and call the vanilla sample as cluster center, because augmentations are produced by using the vanilla sample  $\mathbf{c}_i$  and the augmentation  $\mathbf{x}$  obeys  $\|\mathbf{x} - \mathbf{c}_i\|_2 \leq \epsilon_0$ .

**Definition 1** ( $(\epsilon, \delta)$ -clusterable dataset). *Suppose  $\{(\mathbf{x}_i, \mathbf{y}_i^*)\}_{i=1}^n$  denote the pairs of augmentation and ground-truth label, where augmentation  $\mathbf{x}_i$  generated from the  $t$ -th sample  $\mathbf{c}_t$  obeys  $\|\mathbf{x} - \mathbf{c}_t\|_2 \leq \epsilon$  with a constant  $\epsilon$ , and  $\mathbf{y}_i^* \in \{\gamma_1, \gamma_2, \dots, \gamma_K\}$  of  $\mathbf{x}_i$  is the label of  $\mathbf{c}_t$ . Moreover, samples and its augmentations are normalized, i.e.  $\|\mathbf{c}_i\|_2 = \|\mathbf{x}_i\|_2 = 1$ . Each vanilla sample  $\mathbf{c}_i$  has  $n_i$  augmentations, where  $c_l \frac{n}{K} \leq n_i \leq c_u \frac{n}{K}$  with two constants  $c_l$  and  $c_u$ . Moreover, the classes are separated such that*

$$|\gamma_r - \gamma_s| \geq \delta, \quad \|\mathbf{c}_r - \mathbf{c}_s\|_2 \geq 2\epsilon, \quad (\forall r \neq s),$$

where  $\delta$  is the label separation.

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

**Assumption 1** (Bimodal Jacobian). Let  $\beta \geq \alpha \geq \epsilon > 0$  be scalars. Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}^n$  be a nonlinear mapping and consider a set  $\mathcal{D} \subset \mathbb{R}^p$  containing the initial point  $\mathbf{w}_0$  (i.e.  $\mathbf{w}_0 \in \mathcal{D}$ ). Let  $\mathcal{S}_+ \subset \mathbb{R}^n$  be a subspace and  $\mathcal{S}_-$  be its complement. We say the mapping  $f$  has a Bimodal Jacobian with respect to the complementary subspaces  $\mathcal{S}_+$  and  $\mathcal{S}_-$  as long as the following two assumptions hold for all  $\mathbf{w} \in \mathcal{D}$ .

- **Spectrum over  $\mathcal{S}_+$** : For all  $\mathbf{v} \in \mathcal{S}_+$  with unit Euclidian norm we have

$$\alpha \leq \|\mathcal{J}^T(\mathbf{w})\mathbf{v}\|_{\ell_2} \leq \beta.$$

- **Spectrum over  $\mathcal{S}_-$** : For all  $\mathbf{v} \in \mathcal{S}_-$  with unit Euclidian norm we have

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

We will refer to  $\mathcal{S}_+$  as the signal subspace and  $\mathcal{S}_-$  as the noise subspace.

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

**Assumption 2** (Low-rank Jacobian). Let  $\beta \geq \alpha > 0$  be scalars. Consider a set  $\mathcal{D} \subset \mathbb{R}^p$  containing the initial point  $\mathbf{w}_0$  (i.e.  $\mathbf{w}_0 \in \mathcal{D}$ ). Let  $\mathcal{S}_+ \subset \mathbb{R}^n$  be a subspace and  $\mathcal{S}_-$  be its complement. For all  $\mathbf{w} \in \mathcal{D}$ ,  $\mathbf{v} \in \mathcal{S}_+$  and  $\mathbf{v}' \in \mathcal{S}_-$  with unit Euclidian norm, we have that

$$\alpha \leq \|\mathcal{J}^T(\mathbf{w})\mathbf{v}\|_{\ell_2} \leq \beta \quad \text{and} \quad \|\mathcal{J}^T(\mathbf{w})\mathbf{v}'\|_{\ell_2} = 0.$$

In Theorem 7, we verify that the Jacobian matrix of real datasets indeed have a bimodal structure i.e. there are few large singular values and the remaining singular values are small which further motivate Assumption 2. This is inline with earlier papers which observed that Hessian matrices of deep networks have bimodal spectrum (approximately low-rank) [17] and is related to various results demonstrating that there are flat directions in the loss landscape [18].

Our dataset model in Definition 1 naturally has a low-rank Jacobian when  $\epsilon_0 = 0$  and each augmentation is equal to one of the  $K$  centers (vanilla samples)  $\{\mathbf{c}_\ell\}_{\ell=1}^K$ . In this case, the Jacobian will be at most rank  $K$  since each row will be in the span of  $\{\frac{\partial f(\mathbf{c}_\ell, \mathbf{w})}{\partial \mathbf{w}}\}_{\ell=1}^K$ . The subspace  $\mathcal{S}_+$  is dictated by the membership of each cluster center (vanilla example) as follows: Let  $\Lambda_\ell \subset \{1, \dots, n\}$  be the set of coordinates  $i$  such that  $\mathbf{x}_i = \mathbf{c}_\ell$ . Then, subspace is characterized by  $\mathcal{S}_+ = \{\mathbf{v} \in \mathbb{R}^n \mid \mathbf{v}_{i_1} = \mathbf{v}_{i_2} \text{ for all } i_1, i_2 \in \Lambda_\ell \text{ and } 1 \leq \ell \leq K\}$ . When  $\epsilon_0 > 0$  and the augmentation points of each cluster (vanilla sample) are not the same as the cluster we have the bimodal Jacobian structure of Assumption 1 where over  $\mathcal{S}_-$  the spectral norm is small but nonzero.

**Definition 2** (Support subspace). Let  $\{\mathbf{x}_i\}_{i=1}^n$  be an input dataset generated according to Definition 1. Also let  $\{\tilde{\mathbf{x}}_i\}_{i=1}^n$  be the associated vanilla samples, that is,  $\tilde{\mathbf{x}}_i = \mathbf{c}_\ell$  iff  $\mathbf{x}_i$  is from the  $\ell$ th vanilla sample. We define the support subspace  $\mathcal{S}_+$  as a subspace of dimension  $K$ , dictated by the cluster center membership as follows. Let  $\Lambda_\ell \subset \{1, \dots, n\}$  be the set of coordinates  $i$  such that  $\tilde{\mathbf{x}}_i = \mathbf{c}_\ell$ . Then,  $\mathcal{S}_+$  is characterized by

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

Before we state our general result we need to discuss another assumption and definition.

**Assumption 3** (Smoothness). The Jacobian mapping  $\mathcal{J}(\mathbf{w})$  associated to a nonlinear mapping  $f : \mathbb{R}^p \rightarrow \mathbb{R}^n$  is  $L$ -smooth if for all  $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^p$  we have  $\|\mathcal{J}(\mathbf{w}_2) - \mathcal{J}(\mathbf{w}_1)\| \leq L \|\mathbf{w}_2 - \mathbf{w}_1\|_{\ell_2}$ .

In Theorem 7, we verify this assumption. Note that, if  $\frac{\partial \mathcal{J}(\mathbf{w})}{\partial \mathbf{w}}$  is continuous, the smoothness condition holds over any compact domain (albeit for a possibly large  $L$ ).

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

**Definition 3** (Diffusedness).  $\mathcal{S}_+$  is  $\zeta$  diffused if for any vector  $\mathbf{v} \in \mathcal{S}_+$

$$\|\mathbf{v}\|_\infty \leq \sqrt{\zeta/n} \|\mathbf{v}\|_2,$$

holds for some  $\zeta > 0$ .

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

**Definition 4** (Average Jacobian). *We define the average Jacobian along the path connecting two points  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$  as*

$$\mathcal{J}(\mathbf{y}, \mathbf{x}) := \int_0^1 \mathcal{J}(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x})) d\alpha.$$

**Definition 5** (Neural Net Jacobian). *Given input samples  $(\mathbf{x}_i)_{i=1}^n$ , form the input matrix  $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$ . The Jacobian of our learning problem, i.e.  $\mathbf{x} \mapsto f(\mathbf{W}, \mathbf{x}) = \mathbf{v}^T \phi(\mathbf{W}\mathbf{x})$  and  $\mathcal{L}_t(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_{ti} - f(\mathbf{W}, \mathbf{x}_i))^2$ , at a matrix  $\mathbf{W}$  is denoted by  $\mathcal{J}(\mathbf{W}, \mathbf{X}) \in \mathbb{R}^{n \times kd}$  and is given by*

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

Here  $*$  denotes the Khatri-Rao product.

### C.1.3 Auxiliary Lemmas

**Lemma 3** (Linearization of the residual). *For the general problem (10) in Appendix C.1.1, we define*

$$\mathbf{G}(\mathbf{w}_t) = \mathcal{J}(\mathbf{w}_{t+1}, \mathbf{w}_t)\mathcal{J}(\mathbf{w}_t)^T.$$

where  $\mathcal{J}(\mathbf{w}_t)$  denotes the Jacobian matrix defined in Eqn. (12), and  $\mathcal{J}(\mathbf{w}_{t+1}, \mathbf{w}_t) = \int_0^1 \mathcal{J}(\mathbf{w}_t + \alpha(\mathbf{w}_{t+1} - \mathbf{w}_t)) d\alpha$  denotes the average Jacobian matrix defined in Definition (4). When using the gradient descent iterate  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla \mathcal{L}_t(\mathbf{w}_t)$ , then residuals

$$\mathbf{r}_{t+1} = f(\mathbf{w}_{t+1}) - \bar{\mathbf{y}}^{t+1}, \quad \mathbf{r}_t = f(\mathbf{w}_t) - \bar{\mathbf{y}}^t$$

obey the following equation

$$\mathbf{r}_{t+1} = (\mathbf{I} - \eta \mathbf{G}(\mathbf{w}_t))\mathbf{r}_t + \bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1}.$$

*Proof.* Here we follow [19] to prove our result. Following Definition 4, denoting  $\mathbf{r}_{t+1} = f(\mathbf{w}_{t+1}) - \bar{\mathbf{y}}^{t+1}$  and  $\mathbf{r}_t = f(\mathbf{w}_t) - \bar{\mathbf{y}}^t$ , we find that

$$\begin{aligned} \mathbf{r}_{t+1} &= \mathbf{r}_t - f(\mathbf{w}_t) + f(\mathbf{w}_{t+1}) + \bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1} \\ &\stackrel{\textcircled{1}}{=} \mathbf{r}_t + \mathcal{J}(\mathbf{w}_{t+1}, \mathbf{w}_t)(\mathbf{w}_{t+1} - \mathbf{w}_t) + \bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1} \\ &\stackrel{\textcircled{2}}{=} \mathbf{r}_t - \eta \mathcal{J}(\mathbf{w}_{t+1}, \mathbf{w}_t)\mathcal{J}(\mathbf{w}_t)^T \mathbf{r}_t + \bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1} \\ &= (\mathbf{I} - \eta \mathbf{G}(\mathbf{w}_t))\mathbf{r}_t + \bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1}. \end{aligned}$$

where  $\textcircled{1}$  uses the fact that Jacobian is the derivative of  $f$  and  $\textcircled{2}$  uses the fact that  $\nabla \mathcal{L}_t(\mathbf{w}) = \mathcal{J}(\mathbf{w})^T \mathbf{r}_t$ .  $\square$

Using Assumption 3, one can show that sparse vectors have small projection on  $\mathcal{S}_+$ .

**Lemma 4.** [19] *Suppose Assumption 3 holds. If  $\mathbf{r} \in \mathbb{R}^n$  is a vector with  $s$  nonzero entries, we have that*

$$\|\mathcal{P}_{\mathcal{S}_+}(\mathbf{r})\|_\infty \leq \frac{\zeta \sqrt{s}}{n} \|\mathbf{r}\|_2,$$

where  $\mathcal{P}_{\mathcal{S}_+}(\mathbf{r})$  projects  $\mathbf{r}$  onto the space  $\mathcal{S}_+$ .

**Lemma 5.** *For the general problem (10) in Appendix C.1.1, let  $\mathbf{r}_t = f(\mathbf{w}_t) - \bar{\mathbf{y}}^t$  and  $\hat{\mathbf{r}}_t = \mathcal{P}_{\mathcal{S}_+}(\mathbf{r}_t)$ . Suppose Assumption 2 holds and  $\eta \leq \frac{1}{\beta^2}$ . If  $\|\mathbf{w}_t - \mathbf{w}_0\|_2 + \frac{\|\hat{\mathbf{r}}_t\|_2}{\alpha} \leq \frac{4(1+\psi)\|\mathbf{r}_0\|_2}{\alpha}$ , then*

$$\mathbf{w}_{t+1} \in \mathcal{D} = \left\{ \mathbf{w} \in \mathbb{R}^p \mid \|\mathbf{w} - \mathbf{w}_0\|_2 \leq \frac{4(1+\psi)\|\mathbf{r}_0\|_2}{\alpha} \right\}.$$

*Proof.* Since range space of Jacobian is in  $\mathcal{S}_+$  and  $\eta \leq 1/\beta^2$ , we can easily obtain

$$\begin{aligned}
\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2 &= \eta \|\mathcal{J}^T(\mathbf{w}_t) (f(\mathbf{w}_t) - \bar{\mathbf{y}}^t)\|_2 \\
&\stackrel{\textcircled{1}}{=} \eta \|\mathcal{J}^T(\mathbf{w}_t) (\mathcal{P}_{\mathcal{S}_+}(f(\mathbf{w}_t) - \bar{\mathbf{y}}^t))\|_2 \\
&\stackrel{\textcircled{2}}{=} \eta \|\mathcal{J}^T(\mathbf{w}_t) \hat{\mathbf{r}}_t\|_2 \\
&\stackrel{\textcircled{3}}{\leq} \eta \beta \|\hat{\mathbf{r}}_t\|_2 \\
&\stackrel{\textcircled{4}}{\leq} \frac{\|\hat{\mathbf{r}}_t\|_2}{\beta} \\
&\stackrel{\textcircled{5}}{\leq} \frac{\|\hat{\mathbf{r}}_t\|_2}{\alpha}
\end{aligned}$$

In the above,  $\textcircled{1}$  follows from the fact that row range space of Jacobian is subset of  $\mathcal{S}_+$  via Assumption 2.  $\textcircled{2}$  follows from the definition of  $\hat{\mathbf{r}}_t = \mathcal{P}_{\mathcal{S}_+}(f(\mathbf{w}_t) - \bar{\mathbf{y}}^t)$ .  $\textcircled{3}$  follows from the upper bound on the spectral norm of the Jacobian over  $\mathcal{D}$  per Assumption 2,  $\textcircled{4}$  from the fact that  $\eta \leq \frac{1}{\beta^2}$ ,  $\textcircled{5}$  from  $\alpha \leq \beta$ . The latter combined with the triangular inequality and the assumption

$$\|\mathbf{w}_{t+1} - \mathbf{w}_0\|_2 \leq \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2 + \|\mathbf{w}_0 - \mathbf{w}_t\|_2 \leq \|\mathbf{w}_t - \mathbf{w}_0\|_2 + \frac{\|\hat{\mathbf{r}}_t\|_2}{\alpha} \leq \frac{4(1 + \psi)\|\mathbf{r}_0\|_2}{\alpha},$$

concluding the proof of  $\mathbf{r}_{t+1} \in \mathcal{D}$ .  $\square$

**Lemma 6.** [19] Let  $\mathcal{P}_{\mathcal{S}_+} \in \mathbb{R}^{n \times n}$  be the projection matrix to  $\mathcal{S}_+$  i.e. it is a positive semi-definite matrix whose eigenvectors over  $\mathcal{S}_+$  is 1 and its complement is 0. Let  $\mathbf{r}_t = f(\mathbf{w}_t) - \mathbf{y}_t$ ,  $\hat{\mathbf{r}}_t = \mathcal{P}_{\mathcal{S}_+}(\mathbf{r}_t)$ , and  $\mathbf{G}(\mathbf{w}_t) = \mathcal{J}(\mathbf{w}_{t+1}, \mathbf{w}_t) \mathcal{J}(\mathbf{w}_t)^T$ . Suppose Assumptions 2 and 3 hold, the learning rate  $\eta$  satisfies  $\eta \leq \frac{\alpha}{L\beta\|\mathbf{r}_0\|_2}$ ,  $\|\hat{\mathbf{r}}_t\|_2 \leq \|\mathbf{r}_0\|_2$ , then it holds

$$\beta^2 \mathcal{P}_{\mathcal{S}_+} \succeq \mathbf{G}(\mathbf{w}_t) \succeq \frac{1}{2} \mathcal{J}(\mathbf{w}_t) \mathcal{J}(\mathbf{w}_t)^T \succeq \frac{\alpha^2}{2} \mathcal{P}_{\mathcal{S}_+}.$$

In the above context, we focus on introducing theoretical results for the general problem (10) in Appendix C.1.1. Now we introduce lemmas and theories for our network learning problem, i.e.  $\mathbf{x} \mapsto f(\mathbf{W}, \mathbf{x}) = \mathbf{v}^T \phi(\mathbf{W}\mathbf{x})$  and  $\mathcal{L}_t(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i^t - f(\mathbf{W}, \mathbf{x}_i))^2$  used in our manuscript. Specifically, we introduce some theoretical results in [20] and characterizes three key properties of the neural network Jacobian. These are smoothness, spectral norm, and minimum singular value at initialization which correspond to Lemmas 6.6, 6.7, and 6.8 in that paper.

**Theorem 3** (Jacobian Properties at Cluster Center). [20] Suppose  $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$  be an input dataset satisfying  $\lambda(\mathbf{X}) > 0$ , where  $\lambda(\mathbf{X})$  denotes the smallest eigenvalue of matrix  $\mathbf{X}$ . Suppose  $|\phi'|, |\phi''| \leq \Gamma$  where  $\phi'$  and  $\phi''$  respectively denotes the first and second order derivatives. The Jacobian mapping with respect to the input-to-hidden weights obey the following properties. Let  $\mathcal{J}(\mathbf{W}, \mathbf{X})$  denote the neural net Jacobian defined in Definition 5.

(1) Smoothness is bounded by

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

(2) Top singular value is bounded by

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

(3) Let  $C > 0$  be an absolute constant. As long as

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

At random Gaussian initialization  $\mathbf{W}_0 \sim \mathcal{N}(0, 1)^{k \times d}$ , with probability at least  $1 - 1/K^{100}$ , we have

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

The following theorem states the properties of the Jacobian at a  $(\epsilon_0, \delta)$  clusterable dataset defined in Definition 1. That is,  $(\mathbf{x}_i)_{i=1}^n$  are generated from  $(\mathbf{c}_i)_{i=1}^K$ , and their augmentation distance is at most  $\epsilon_0$  and label separation is at least  $\delta$ .

**Theorem 4** (Jacobian Properties at Cluster Center). [19] *Let input samples  $(\mathbf{x}_i)_{i=1}^n$  be generated according to  $(\epsilon_0, \delta)$  clusterable dataset model of Definition 1. Define  $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]^T$  and  $\mathbf{C} = [\mathbf{c}_1 \dots \mathbf{c}_k]^T$ . Let  $\mathcal{S}_+$  be the support space and  $(\tilde{\mathbf{x}}_i)_{i=1}^n$  be the associated clean dataset as described by Definition 2. Set  $\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1 \dots \tilde{\mathbf{x}}_n]^T$ . Assume  $|\phi'|, |\phi''| \leq \Gamma$  and  $\lambda(\mathbf{C}) > 0$ . Let  $\mathcal{J}(\mathbf{W}, \mathbf{X})$  denote the neural net Jacobian defined in Definition 5. The Jacobian mapping at  $\tilde{\mathbf{X}}$  with respect to the input-to-hidden weights obey the following properties.*

(1) *Smoothness is bounded by*

$$\left\| \mathcal{J}(\tilde{\mathbf{W}}, \tilde{\mathbf{X}}) - \mathcal{J}(\mathbf{W}, \tilde{\mathbf{X}}) \right\| \leq \Gamma \sqrt{\frac{c_{up}n}{kK}} \|\mathbf{C}\| \left\| \tilde{\mathbf{W}} - \mathbf{W} \right\|_F \quad \text{for all } \tilde{\mathbf{W}}, \mathbf{W} \in \mathbb{R}^{k \times d}.$$

(2) *Top singular value is bounded by*

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

(3) *As long as*

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

*At random Gaussian initialization  $\mathbf{W}_0 \sim \mathcal{N}(0, 1)^{k \times d}$ , with probability at least  $1 - 1/K^{100}$ , we have*

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

(4) *The range space obeys  $\text{range}(\mathcal{J}(\mathbf{W}_0, \tilde{\mathbf{X}})) \subset \mathcal{S}_+$  where  $\mathcal{S}_+$  is given by Definition 2.*

**Lemma 7** (Upper bound on initial misfit). [19] *Consider a one-hidden layer neural network model of the form  $\mathbf{x} \mapsto \mathbf{v}^T \phi(\mathbf{W}\mathbf{x})$  where the activation  $\phi$  has bounded derivatives obeying  $|\phi(0)|, |\phi'(z)| \leq \Gamma$ . Suppose entries of  $\mathbf{v} \in \mathbb{R}^k$  are half  $1/\sqrt{k}$  and half  $-1/\sqrt{k}$  so that  $\|\mathbf{v}\|_2 = 1$ . Also assume we have  $n$  data points  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$  with unit euclidean norm ( $\|\mathbf{x}_i\|_2 = 1$ ) aggregated as rows of a matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and the corresponding labels given by  $\mathbf{y} \in \mathbb{R}^n$  generated according to  $(\rho, \epsilon = 0, \delta)$  noisy dataset (Definition 1). Then for  $\mathbf{W}_0 \in \mathbb{R}^{k \times d}$  with i.i.d.  $\mathcal{N}(0, 1)$  entries*

$$\|\mathbf{v}^T \phi(\mathbf{W}_0 \mathbf{X}^T) - \mathbf{y}\|_2 \leq \mathcal{O} \left( \Gamma \sqrt{n \log K} \right),$$

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

Then we introduce a lemma regarding the projection of label noise on the vanilla sample (cluster) induced subspace. Since augmentations are produced by using the vanilla sample  $\mathbf{c}_i$  and the augmentation  $\mathbf{x}$  obeys  $\|\mathbf{x} - \mathbf{c}_i\|_2 \leq \epsilon_0$ . So in this sense, we sometimes call the vanilla sample and its augmentations as cluster, and call the vanilla sample as cluster center.

**Lemma 8.** [19] *Let  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  be an  $(\rho, \epsilon = 0, \delta)$  clusterable noisy dataset as described in Definition 1. Let  $\{\mathbf{y}_i^*\}_{i=1}^n$  be the corresponding ground truth labels. Let  $\mathcal{J}(\mathbf{W}, \mathbf{C})$  be the Jacobian at the cluster center matrix which is rank  $K$  and  $\mathcal{S}_+$  be its column space. Then, the difference between noiseless and noisy labels satisfy the bound*

$$\|\mathcal{P}_{\mathcal{S}_+}(\mathbf{y} - \mathbf{y}^*)\|_{\infty} \leq 2\rho.$$

**Theorem 5.** [19] *Assume  $|\phi'|, |\phi''| \leq \Gamma$  and  $k \gtrsim d$ . Suppose  $\mathbf{W}_0 \sim \mathcal{N}(0, 1)$ . Let  $\mathbf{c}_1, \dots, \mathbf{c}_K$  be cluster centers. Then, with probability at least  $1 - 2e^{-(k+d)} - Ke^{-100d}$  over  $\mathbf{W}_0$ , any matrix  $\mathbf{W}$  satisfying  $\|\mathbf{W} - \mathbf{W}_0\|_F \lesssim \sqrt{k}$  satisfies the following. For all  $1 \leq i \leq K$ ,*

$$\sup_{\|\mathbf{x} - \mathbf{c}_i\|_2, \|\tilde{\mathbf{x}} - \mathbf{c}_i\|_2 \leq \epsilon} |f(\mathbf{W}, \mathbf{x}) - f(\mathbf{W}, \tilde{\mathbf{x}})| \leq C\Gamma\epsilon(\|\mathbf{W} - \mathbf{W}_0\| + \sqrt{d}).$$

**Lemma 9** (Perturbed Jacobian Distance). [19] Let  $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]^T$  be the input matrix obtained from Definition 1. Let  $\widetilde{\mathbf{X}}$  be the noiseless inputs where  $\widetilde{\mathbf{x}}_i$  is the cluster center corresponding to  $\mathbf{x}_i$ . Let  $\mathcal{J}(\mathbf{W}, \mathbf{X})$  denote the neural net Jacobian defined in Definition 5 and define  $\mathcal{J}(\mathbf{W}_1, \mathbf{W}_2, \mathbf{X}) = \int_0^1 \mathcal{J}(\alpha \mathbf{W}_1 + (1 - \alpha) \mathbf{W}_2, \mathbf{X}) d\alpha$ . Given weight matrices  $\mathbf{W}_1, \mathbf{W}_2, \widetilde{\mathbf{W}}_1, \widetilde{\mathbf{W}}_2$ , we have that

$$\|\mathcal{J}(\mathbf{W}, \mathbf{X}) - \mathcal{J}(\widetilde{\mathbf{W}}, \widetilde{\mathbf{X}})\| \leq \Gamma \sqrt{n} \left( \frac{\|\widetilde{\mathbf{W}} - \mathbf{W}\|_F}{\sqrt{k}} + \varepsilon \right).$$

and

$$\|\mathcal{J}(\mathbf{W}_1, \mathbf{W}_2, \mathbf{X}) - \mathcal{J}(\widetilde{\mathbf{W}}_1, \widetilde{\mathbf{W}}_2, \widetilde{\mathbf{X}})\| \leq \Gamma \sqrt{n} \left( \frac{\|\widetilde{\mathbf{W}}_1 - \mathbf{W}_1\|_F + \|\widetilde{\mathbf{W}}_2 - \mathbf{W}_2\|_F}{2\sqrt{k}} + \varepsilon \right).$$

## C.2 Proof of Theorem 2

The subsection has four parts. In the first part, we introduce the proof roadmap in Appendix C.2.1. Then in the second part, we present several auxiliary theories in Appendix C.2.2. Next, we prove our Theorem 2 in Appendix C.2.3. Finally, we present all proof details of auxiliary theories in Appendix C.2.2.

### C.2.1 Proof roadmap

Before proving Theorem 2, we first briefly introduce our main idea. In the **first step**, we analyze the general model introduced in Appendix C.1.1. For the solution  $\mathbf{w}_t$  at the  $t$ -th iteration, Theorem 6 proves that (1) the distance of  $\|\mathbf{w}_t - \mathbf{w}_0\|_2$  can be upper bounded; (2) both residual  $\|\mathcal{P}_{S_+}(f(\mathbf{w}_t) - \widetilde{\mathbf{y}}^t)\|_2$  and  $\|f(\mathbf{w}_t) - \mathbf{y}^*\|_\infty$  can be upper bound. Result (1) means that the gradient descent algorithm gives solutions in a ball around the initialization  $\mathbf{w}_0$ , and helps us verify our assumptions, e.g. Assumptions 3 and 2 and upper bound some variables in our analysis. Results (2) directly bound the label estimation error which plays key role in subsequent analysis.

In the **second step**, we prove Theorem 7 for the perfectly clustered data ( $\epsilon_0 = 0$ ) by using Theorem 6. We consider  $\epsilon_0 \rightarrow 0$  which means that the input data set is perfectly clean. In this setting, let  $\widetilde{\mathbf{X}} = [\widetilde{\mathbf{x}}_1, \dots, \widetilde{\mathbf{x}}_n]$  be the clean input sample matrix obtained by mapping  $\mathbf{x}_i$  to its associated cluster center, i.e.  $\widetilde{\mathbf{x}}_i = \mathbf{c}_\ell$  if  $\mathbf{x}_i$  belongs to the  $\ell$ -th cluster. In this way, we update network parameter  $\widetilde{\mathbf{W}}_t$  as follows:

$$\widetilde{\mathbf{W}}_{t+1} = \widetilde{\mathbf{W}}_t - \nabla \widetilde{\mathcal{L}}_t(\widetilde{\mathbf{W}}_t) \quad \text{where} \quad \widetilde{\mathcal{L}}_t(\widetilde{\mathbf{W}}) = \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_{ti} - f(\widetilde{\mathbf{W}}, \widetilde{\mathbf{x}}_i))^2$$

Theorem 7 shows that for neural networks, our method still can upper bound the distance  $\|\widetilde{\mathbf{W}}_t - \widetilde{\mathbf{W}}_0\|_F$  and the residuals  $\|f(\widetilde{\mathbf{W}}_t) - \widetilde{\mathbf{y}}\|_\infty$  if the network, learning rate, the weight  $\alpha_i$  for refining label satisfy certain conditions.

In the **third step**, we consider the realistic setting, where we update the parameters on the corrupted data  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$  as follows:

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \nabla \mathcal{L}_t(\mathbf{W}_t) \quad \text{where} \quad \mathcal{L}_t(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_{ti} - f(\mathbf{W}, \mathbf{x}_i))^2. \quad (13)$$

Then to upper bound  $\|f(\mathbf{W}_t) - \widetilde{\mathbf{y}}\|_\infty$  which measures the error between the predicted label  $f(\mathbf{W}_t)$  and the ground truth label  $\widetilde{\mathbf{y}}$ , we upper bound  $\|f(\mathbf{W}_t, \mathbf{X}) - f(\widetilde{\mathbf{W}}_t, \widetilde{\mathbf{X}})\|_2$  and  $\|\mathbf{W}_t - \widetilde{\mathbf{W}}_t\|_F$ . These results are formally stated in Theorem 8.

In the **fourth step**, we combine the above results together. Specifically, Theorem 7 upper bounds the residuals  $\|f(\widetilde{\mathbf{W}}_t, \widetilde{\mathbf{X}}) - \widetilde{\mathbf{y}}\|_\infty$  and Theorem 8 upper bounds  $\|f(\mathbf{W}_t, \mathbf{X}) - f(\widetilde{\mathbf{W}}_t, \widetilde{\mathbf{X}})\|_2$ . So combining these two results and other results in Theorem 7 & 8, we can upper bound  $\|f(\mathbf{W}_t, \mathbf{X}) - \widetilde{\mathbf{y}}\|_\infty$  which is our desired results. At the same time, by using similar method, we can also bound the label estimation error by our self-labeling refinery, since  $\|\widetilde{\mathbf{y}}^t - \mathbf{y}^*\|_2 = \|(1 - \alpha_t)\mathbf{y} + \alpha_t f(\mathbf{w}) - \mathbf{y}^*\|_2 \leq (1 - \alpha_t)\|\mathbf{y} - \mathbf{y}^*\|_2 + \alpha_t \|f(\mathbf{w}) - \mathbf{y}^*\|_2$ . The term  $\|\mathbf{y} - \mathbf{y}^*\|_2$  denotes the initial label error and can be bounded by a factor related to  $\rho$ , while the second term is well upper bounded by the above results.

It should be note that our proof framework follows the recent works [21, 19] which shows that gradient descent is robust to label corruptions. The main difference is that this work uses the label estimation  $\tilde{\mathbf{y}}^t = \alpha_t \mathbf{y} + (1 - \alpha_t) f(\mathbf{w})$  and minimizes the squared loss, while both works [21, 19] use the corrupted label  $\mathbf{y}$  and then minimize the squared loss. By comparison, our method is much more complicated and gives different proofs.

### C.2.2 Auxiliary Theories

The following theorem is to analyze the general model introduced in Appendix C.1.1. It guarantees that the estimated label by our method is close to the ground truth label when the Jacobian mapping is exactly low-rank. By using this results, one can obtain Theorem 7 for the perfectly clustered data ( $\epsilon_0 = 0$ ) which will be stated later.

**Theorem 6** (Gradient descent with label corruption). *Consider a nonlinear least squares problem of the form  $\mathcal{L}_t(\mathbf{w}) = \frac{1}{2} \|f(\mathbf{w}) - \tilde{\mathbf{y}}^t\|_{\ell_2}^2$  with the nonlinear mapping  $f : \mathbb{R}^p \rightarrow \mathbb{R}^n$  obeying assumptions 2 and 3 over a unit Euclidian ball of radius  $\frac{4(1+\psi_1)\|f(\mathbf{w}_0) - \mathbf{y}\|_2}{\alpha}$  around an initial point  $\mathbf{w}_0$  and  $\mathbf{y} = [y_1 \dots y_n] \in \mathbb{R}^n$  denoting the corrupted labels. We also assume  $\alpha_t \geq 1 - \frac{\alpha^2}{4\beta^2}$  and  $2\sqrt{n} \lim_{t \rightarrow +\infty} \sum_{t=0}^t |\alpha_t - \alpha_{t+1}| \leq \psi_1 \|f(\mathbf{w}_0) - \tilde{\mathbf{y}}^0\|_2$ . Also let  $\mathbf{y}^* = [\mathbf{y}_1^* \dots \mathbf{y}_n^*] \in \mathbb{R}^n$  denote the ground truth labels and  $\mathbf{e} = \mathbf{y} - \mathbf{y}^*$  the corruption. Furthermore, suppose the initial residual  $f(\mathbf{w}_0) - \tilde{\mathbf{y}}$  with respect to the uncorrupted labels obey  $f(\mathbf{w}_0) - \mathbf{y}^* \in \mathcal{S}_+$ . Then, running gradient descent updates of the from (11) with a learning rate  $\eta \leq \frac{1}{2\beta^2} \min\left(1, \frac{\alpha\beta}{L\|f(\mathbf{w}_0) - \tilde{\mathbf{y}}^0\|_2}\right)$ , all iterates obey*

$$\|\mathbf{w}_t - \mathbf{w}_0\|_2 \leq \frac{4\|\mathbf{r}_0\|_2}{\alpha} + 2\sqrt{n} \lim_{t \rightarrow +\infty} \sum_{t=0}^t |\alpha_t - \alpha_{t+1}| \leq \frac{4(1+\psi)\|f(\mathbf{w}_0) - \tilde{\mathbf{y}}^0\|_2}{\alpha}.$$

and

$$\|\hat{\mathbf{r}}_t\|_2^2 \leq \left(1 - \frac{\eta\alpha^2}{4}\right)^t \|\hat{\mathbf{r}}_0\|_2^2 + 2\sqrt{n} \sum_{i=0}^{t-1} \left(1 - \frac{\eta\alpha^2}{4}\right)^{t-i} |\alpha_i - \alpha_{i+1}|,$$

where  $\mathbf{r}_t = f(\mathbf{w}_t) - \tilde{\mathbf{y}}^t$  and let  $\mathbf{r}_0 = f(\mathbf{w}_0) - \tilde{\mathbf{y}}^0$  be the initial residual, and  $\hat{\mathbf{r}}_t = \mathcal{P}_{\mathcal{S}_+}(\mathbf{r}_t)$ . Furthermore, assume  $\nu > 0$  is a precision level obeying  $\nu \geq \|\mathcal{P}_{\mathcal{S}_+}(\mathbf{e})\|_\infty$ . Then, after  $t \geq \frac{5}{\eta\alpha^2} \log\left(\frac{\|f(\mathbf{w}_0) - \tilde{\mathbf{y}}^0\|_2}{(1-\alpha_{\max})^\nu}\right)$  iterations where  $\alpha_{\max} = \max_t \alpha_t$ ,  $\mathbf{w}_t$  achieves the following error bound with respect to the true labels

$$\|f(\mathbf{w}_t) - \mathbf{y}^*\|_\infty \leq 2\nu + \frac{2\sqrt{n}}{1 - \alpha_t} \sum_{i=0}^{t-1} \left(1 - \frac{\eta\alpha^2}{4}\right)^{t-i} |\alpha_i - \alpha_{i+1}|.$$

Furthermore, if  $\mathbf{e}$  has at most  $s$  nonzeros and  $\mathcal{S}_+$  is  $\zeta$  diffused per Definition 3, then using  $\nu = \|\mathcal{P}_{\mathcal{S}_+}(\mathbf{e})\|_\infty$

$$\begin{aligned} \|f(\mathbf{w}_t) - \mathbf{y}^*\|_\infty &\leq 2\|\mathcal{P}_{\mathcal{S}_+}(\mathbf{e})\|_\infty + \frac{2\sqrt{n}}{1 - \alpha_t} \sum_{i=0}^{t-1} \left(1 - \frac{\eta\alpha^2}{4}\right)^{t-i} |\alpha_i - \alpha_{i+1}| \\ &\leq \frac{\zeta\sqrt{s}}{n} \|\mathbf{e}\|_2 + \frac{2\sqrt{n}}{1 - \alpha_t} \sum_{i=0}^{t-1} \left(1 - \frac{\eta\alpha^2}{4}\right)^{t-i} |\alpha_i - \alpha_{i+1}|, \end{aligned}$$

where  $\mathcal{P}_{\mathcal{S}_+}(\mathbf{e})$  denotes projection of  $\mathbf{e}$  on  $\mathcal{S}_+$ .

See its proof in Appendix C.2.5. This result shows that when the Jacobian of the nonlinear mapping is low-rank, our method enjoys two good properties.

For the solution  $\mathbf{w}_t$  at the  $t$ -th iteration, (1) the distance of  $\|\mathbf{w}_t - \mathbf{w}_0\|_2$  can be upper bounded; (2) both residual  $\|\mathcal{P}_{\mathcal{S}_+}(f(\mathbf{w}_t) - \mathbf{y}_t)\|_2$  and  $\|f(\mathbf{w}_t) - \tilde{\mathbf{y}}\|_\infty$  can be upper bound. Result (1) means that the gradient descent algorithm gives solutions in a ball around the initialization  $\mathbf{w}_0$ , and helps us verify our assumptions, e.g. Assumptions 3 and 2 and upper bound some variables in our analysis. Results (2) directly bound the label estimation error which plays key role in subsequent analysis. This theorem is the key result that allows us to prove Theorem 7 when the data points are perfectly

clustered ( $\epsilon_0 = 0$ ). Furthermore, this theorem when combined with a perturbation analysis allows us to deal with data that is not perfectly clustered ( $\epsilon_0 > 0$ ) and to conclude the recovery ability of our method (Theorem 2).

When  $\epsilon_0 \rightarrow 0$  which means that the input data set is perfectly clustered, our method can be expected to exactly recover the ground truth label by using neural networks.

**Theorem 7** (Training with perfectly clustered data). *Consider the setting and assumptions of Theorem 6 with  $\epsilon_0 = 0$ . Starting from an initial weight matrix  $\mathbf{w}_0$  selected at random with i.i.d.  $\mathcal{N}(0, 1)$  entries we run gradient descent updates of the form  $\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \nabla \mathcal{L}_t(\mathbf{W}_t)$  on the least-squares loss in the manuscript with step size  $\eta \leq \frac{K}{2c_{up}n\Gamma^2\|\mathbf{C}\|^2}$ . Furthermore, assume the number of hidden nodes obey*

$$k \geq C(1 + \psi_1)^2 \Gamma^4 \frac{K \log(K) \|\mathbf{C}\|^2}{\lambda(\mathbf{C})^2},$$

with  $\lambda(\mathbf{C})$  is the minimum eigenvalue of  $\Sigma(\mathbf{C})$  in Assumption 2. Then, with probability at least  $1 - 2/K^{100}$  over randomly initialized  $\mathbf{W}_0 \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$ , the iterates  $\mathbf{W}_t$  obey the following properties.

(1) The distance to initial point  $\mathcal{W}_0$  is upper bounded by

$$\|\mathbf{W}_t - \mathbf{W}_0\|_F \leq c\Gamma \sqrt{\frac{K \log K}{\lambda(\mathbf{C})}}.$$

(2) After  $t \geq t_0 := \frac{cK}{\eta n \lambda(\mathbf{C})} \log\left(\frac{\Gamma \sqrt{n \log K}}{(1 - \alpha_{max})\rho}\right)$  iterations where  $\alpha_{max} = \max_{0 \leq t \leq t_0} \alpha_t$ , the entrywise predictions of the learned network with respect to the ground truth labels  $\{\mathbf{y}_i^*\}_{i=1}^n$  satisfy

$$|f(\mathbf{W}_t, \mathbf{x}_i) - \mathbf{y}_i^*| \leq 4\rho,$$

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

$$\arg \min_{i: 1 \leq i \leq K} |f(\mathbf{W}_t, \mathbf{x}_i) - \gamma_i| = \mathbf{y}_i^* \quad \text{for } i = 1, 2, \dots, n. \quad (14)$$

See its proof in Appendix C.2.6. This result shows that in the limit  $\epsilon_0 \rightarrow 0$  where the data points are perfectly clustered, if the width of network and the iterations satisfy  $k \geq C(1 + \psi_1)^2 \Gamma^4 \frac{K \log(K) \|\mathbf{C}\|^2}{\lambda(\mathbf{C})^2}$

and  $t \geq t_0 := \frac{cK}{\eta n \lambda(\mathbf{C})} \log\left(\frac{\Gamma \sqrt{n \log K}}{\alpha \rho}\right)$ , then our method can exactly recover the ground truth label. This result can be interpreted as ensuring that the network has enough capacity to fit the cluster centers  $\{\mathbf{c}_\ell\}_{\ell=1}^K$  and the associated true labels.

Then we consider the perturbed data  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$  instead of the perfectly clustered data  $\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n]$  obtained by mapping  $\mathbf{x}_i$  to its associated cluster center, i.e.  $\tilde{\mathbf{x}}_i = \mathbf{c}_\ell$  if  $\mathbf{x}_i$  belongs to the  $\ell$ -th cluster. In Theorem 8, we upper bound the parameter distance and output distance under the two kinds of data  $\mathbf{X}$  and  $\tilde{\mathbf{X}}$ .

**Theorem 8** (Robustness of gradient path to perturbation). *Generate samples  $(\mathbf{x}_i, \mathbf{y}_i)_{i=1}^n$  according to  $(\rho, \epsilon, \delta)$  corrupted dataset and form the concatenated input/labels  $\mathbf{X} \in \mathbb{R}^{d \times n}$ ,  $\mathbf{y} \in \mathbb{R}^n$ . Let  $\tilde{\mathbf{X}}$  be the clean input sample matrix obtained by mapping  $\mathbf{x}_i$  to its associated cluster center. Set learning rate  $\eta \leq \frac{K}{2c_{up}n\Gamma^2\|\mathbf{C}\|^2}$  and maximum iterations  $t_0$  satisfying*

$$\eta t_0 = C_1 \frac{K}{n \lambda(\mathbf{C})} \log\left(\frac{\Gamma \sqrt{n \log K}}{\rho}\right).$$

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

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

where  $\alpha_{\max} = \max_{1 \leq t \leq t_0} \alpha_t$ . Assume  $2\sqrt{n} \sum_{i=0}^{t-1} \left(1 - \frac{\eta\alpha^2}{4}\right)^{t-i} |\alpha_i - \alpha_{i+1}| \leq \psi_2 \|\mathbf{r}_0\|_2^2$  and  $2\sqrt{n} \sum_{i=0}^{t-1} |\alpha_i - \alpha_{i+1}| \leq \psi_1 \|\mathbf{r}_0\|_2$ . Set  $\mathbf{W}_0 \sim \mathcal{N}(0, 1)$ . Starting from  $\mathbf{W}_0 = \widetilde{\mathbf{W}}_0$  consider the gradient descent iterations over the losses

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \nabla \mathcal{L}_t(\mathbf{W}_t) \quad \text{where} \quad \mathcal{L}_t(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_{ti} - f(\mathbf{W}, \mathbf{x}_i))^2 \quad (15)$$

$$\widetilde{\mathbf{W}}_{t+1} = \widetilde{\mathbf{W}}_t - \nabla \widetilde{\mathcal{L}}_t(\widetilde{\mathbf{W}}_t) \quad \text{where} \quad \widetilde{\mathcal{L}}_t(\widetilde{\mathbf{W}}) = \frac{1}{2} \sum_{i=1}^n (\widetilde{\mathbf{y}}_{ti} - f(\widetilde{\mathbf{W}}, \widetilde{\mathbf{x}}_i))^2 \quad (16)$$

Then, for all gradient descent iterations satisfying  $t \leq t_0$ , we have that

$$\|f(\mathbf{W}_t, \mathbf{X}) - f(\widetilde{\mathbf{W}}_t, \widetilde{\mathbf{X}})\|_2 \leq c_0 \psi' t \eta \varepsilon \Gamma^3 n^{3/2} \sqrt{\log K},$$

and

$$\|\mathbf{W}_t - \widetilde{\mathbf{W}}_t\|_F \leq \mathcal{O} \left( t \psi' \eta \varepsilon \frac{\Gamma^4 K n}{\lambda(\mathbf{C})} \log \left( \frac{\Gamma \sqrt{n \log K}}{\rho} \right)^2 \right).$$

where  $\psi' = 1 + \frac{\psi_1}{2} + \sqrt{\psi_2}$ .

See its proof in Appendix C.2.7. Theorem 2 is obtained by combining the above results together.

### C.2.3 Proof of Theorem 2

*Proof of Theorem 2.* Here we prove our results by three steps. In these steps, each step proves one of the three results in our theory. To begin with, we consider two parameter update settings with initialization as  $\mathbf{W}_0$ :

$$\begin{aligned} \widetilde{\mathbf{W}}_{t+1} &= \widetilde{\mathbf{W}}_t - \nabla \widetilde{\mathcal{L}}_t(\widetilde{\mathbf{W}}_t) \quad \text{where} \quad \widetilde{\mathcal{L}}_t(\widetilde{\mathbf{W}}) = \frac{1}{2} \sum_{i=1}^n (\widetilde{\mathbf{y}}_i^t - f(\widetilde{\mathbf{W}}, \widetilde{\mathbf{x}}_i))^2, \\ \mathbf{W}_{t+1} &= \mathbf{W}_t - \eta \nabla \mathcal{L}_t(\mathbf{W}_t) \quad \text{where} \quad \mathcal{L}_t(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^n (\widetilde{\mathbf{y}}_i^t - f(\mathbf{W}, \mathbf{x}_i))^2, \end{aligned}$$

where  $\widetilde{\mathbf{y}}_i^t = (1 - \alpha_t) \mathbf{y} + \alpha_t f(\widetilde{\mathbf{W}}_t, \widetilde{\mathbf{x}}_i)$ ,  $\widetilde{\mathbf{y}}_i^t = (1 - \alpha_t) \mathbf{y} + \alpha_t f(\mathbf{W}_t, \mathbf{x}_i)$ ,  $\widetilde{\mathbf{X}} = [\widetilde{\mathbf{x}}_1, \dots, \widetilde{\mathbf{x}}_n]$  denotes the clean input sample matrix obtained by mapping  $\mathbf{x}_i$  to its associated cluster center, i.e.  $\widetilde{\mathbf{x}}_i = \mathbf{c}_\ell$  if  $\mathbf{x}_i$  belongs to the  $\ell$ -th cluster, and  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$  denotes corrupted data matrix. Denote the prediction residual vectors of the noiseless and original problems with respect true ground truth labels  $\mathbf{y}^*$  by  $\widetilde{\mathbf{r}}_t = f(\widetilde{\mathbf{W}}_t, \widetilde{\mathbf{X}}) - \mathbf{y}^*$  and  $\mathbf{r}_t = f(\mathbf{W}_t, \mathbf{X}) - \mathbf{y}^*$  respectively.

Theorem 7 shows that if number of iterations  $t$  and network width receptively satisfy  $t \geq t_0 := \frac{cK}{\eta n \lambda(\mathbf{C})} \log \left( \frac{\Gamma \sqrt{n \log K}}{\alpha \rho} \right)$  and  $k \geq C(1 + \psi_1)^2 \Gamma^4 \frac{K \log(K) \|\mathbf{C}\|}{\lambda(\mathbf{C})^2}$ , then it holds

$$\|\widetilde{\mathbf{r}}_t\|_\infty = \|f(\widetilde{\mathbf{W}}_t, \widetilde{\mathbf{X}}) - \mathbf{y}^*\|_\infty \leq 4\rho \quad \text{and} \quad \|\widetilde{\mathbf{W}}_t - \mathbf{W}_0\|_F \leq c\Gamma \sqrt{\frac{K \log K}{\lambda(\mathbf{C})}}.$$

Meanwhile, Theorem 8 proves that if  $\varepsilon \leq \mathcal{O} \left( \frac{\lambda(\mathbf{C})}{\Gamma^2 K \log \left( \frac{\Gamma \sqrt{n \log K}}{\rho} \right)} \right)$  and  $k \geq \mathcal{O} \left( \Gamma^{10} \frac{K^2 \|\mathbf{C}\|^4}{\alpha_{\max}^2 \lambda(\mathbf{C})^4} \log \left( \frac{\Gamma \sqrt{n \log K}}{\rho} \right)^6 \right)$ , then it holds

$$\|\widetilde{\mathbf{r}}_t - \mathbf{r}_t\|_2 \leq c\varepsilon \frac{\psi' K}{n \lambda(\mathbf{C})} \log \left( \frac{\Gamma \sqrt{n \log K}}{\rho} \right) \Gamma^3 n^{3/2} \sqrt{\log K} = c \frac{\psi' \varepsilon \Gamma^3 K \sqrt{n \log K}}{\lambda(\mathbf{C})} \log \left( \frac{\Gamma \sqrt{n \log K}}{\rho} \right)$$

and

$$\|\mathbf{W}_t - \widetilde{\mathbf{W}}_t\|_F \leq \mathcal{O} \left( t \psi' \eta \varepsilon \frac{\Gamma^4 K n}{\lambda(\mathbf{C})} \log \left( \frac{\Gamma \sqrt{n \log K}}{\rho} \right)^2 \right).$$

where  $\psi' = 1 + \frac{\psi_1}{2} + \sqrt{\psi_2}$ .

**Step 1.** By using the above two results, we have

$$\frac{\|f(\mathbf{W}_t, \mathbf{X}) - \tilde{\mathbf{y}}\|_2}{\sqrt{n}} = \frac{1}{\sqrt{n}} (\|\tilde{\mathbf{r}}_t\|_2 + \|\mathbf{r}_t - \tilde{\mathbf{r}}_t\|_2) \leq 4\rho + c \frac{\varepsilon \psi' \Gamma^3 K \sqrt{\log K}}{\lambda(\mathbf{C})} \log\left(\frac{\Gamma \sqrt{n \log K}}{\rho}\right).$$

Moreover, we can also upper bound

$$\begin{aligned} \frac{\|\tilde{\mathbf{y}}^t - \mathbf{y}^*\|_2}{\sqrt{n}} &\leq \frac{(1 - \alpha_t) \|\mathbf{y} - \mathbf{y}^*\|_2}{\sqrt{n}} + \frac{\alpha_t \|f(\mathbf{W}_t, \mathbf{X}) - \mathbf{y}^*\|_2}{\sqrt{n}} \\ &= \frac{(1 - \alpha_t) \|\mathbf{y} - \mathbf{y}^*\|_2}{\sqrt{n}} + 4\alpha_t \rho + c\alpha_t \frac{\varepsilon \psi' \Gamma^3 K \sqrt{\log K}}{\lambda(\mathbf{C})} \log\left(\frac{\Gamma \sqrt{n \log K}}{\rho}\right). \end{aligned}$$

**Step 2.** Now we consider what cases that our method can exactly recover the ground truth label. Assume an input  $\mathbf{x}$  is within  $\varepsilon$ -neighborhood of one of the cluster centers  $\mathbf{c} \in (\mathbf{c}_\ell)_{\ell=1}^K$ . Then we try to upper bound  $|f(\mathbf{W}_t, \mathbf{x}) - f(\tilde{\mathbf{W}}_t, \mathbf{c})|$  where  $f(\tilde{\mathbf{W}}_t, \mathbf{c})$  corresponds to  $f(\tilde{\mathbf{W}}_t, \tilde{\mathbf{x}})$ . To begin with, we have

$$|f(\mathbf{W}_t, \mathbf{x}) - f(\tilde{\mathbf{W}}_t, \mathbf{c})| \leq |f(\mathbf{W}_t, \mathbf{x}) - f(\tilde{\mathbf{W}}_t, \mathbf{x})| + |f(\tilde{\mathbf{W}}_t, \mathbf{x}) - f(\tilde{\mathbf{W}}_t, \mathbf{c})|$$

We upper bound the first term as follows:

$$\begin{aligned} |f(\mathbf{W}_t, \mathbf{x}) - f(\tilde{\mathbf{W}}_t, \mathbf{x})| &= |\mathbf{v}^T \phi(\mathbf{W}_t \mathbf{x}) - \mathbf{v}^T \phi(\tilde{\mathbf{W}}_t \mathbf{x})| \leq \|\mathbf{v}\|_2 \|\phi(\mathbf{W}_t \mathbf{x}) - \phi(\tilde{\mathbf{W}}_t \mathbf{x})\|_2 \\ &\leq \Gamma \|\mathbf{W}_t - \tilde{\mathbf{W}}_t\|_F \\ &\leq \mathcal{O}\left(\varepsilon \psi' \frac{\Gamma^5 K^2}{\lambda(\mathbf{C})^2} \log\left(\frac{\Gamma \sqrt{n \log K}}{\rho}\right)^3\right) \end{aligned}$$

where we use the results  $\|\mathbf{W}_t - \tilde{\mathbf{W}}_t\|_F \leq \mathcal{O}\left(t \psi' \eta \varepsilon \frac{\Gamma^4 K n}{\lambda(\mathbf{C})} \log\left(\frac{\Gamma \sqrt{n \log K}}{\rho}\right)^2\right)$  with  $\psi' = 1 + \frac{\psi_1}{2} + \sqrt{\psi_2}$  in Theorem 8, and  $t = t_0$ . Next, we need to bound

$$|f(\tilde{\mathbf{W}}_t, \mathbf{x}) - f(\tilde{\mathbf{W}}_t, \mathbf{c})| \leq |\mathbf{v}^T \phi(\tilde{\mathbf{W}}_t \mathbf{x}) - \mathbf{v}^T \phi(\tilde{\mathbf{W}}_t \mathbf{c})|.$$

On the other hand, we have  $\|\tilde{\mathbf{W}}_t - \mathbf{W}_0\|_F \leq \mathcal{O}\left(\Gamma \sqrt{\frac{K \log K}{\lambda(\mathbf{C})}}\right)$  in Theorem 7,  $\|\mathbf{x} - \mathbf{c}\|_2 \leq \varepsilon$  and  $\mathbf{W}_0 \sim \mathcal{N}(0, \mathbf{I})$  in assumption. Moreover, using by assumption we have

$$k \geq \mathcal{O}\left(\|\tilde{\mathbf{W}}_t - \mathbf{W}_0\|_F^2\right) = \mathcal{O}\left(\Gamma^2 \frac{K \log K}{\lambda(\mathbf{C})}\right).$$

By using the above results, Theorem 5 guarantees that with probability at  $1 - K \exp(-100d)$ , for all inputs  $\mathbf{x}$  lying  $\varepsilon$  neighborhood of cluster centers, it holds that

$$|f(\mathbf{W}_t, \mathbf{x}) - f(\tilde{\mathbf{W}}_t, \mathbf{c})| \leq C' \Gamma \varepsilon (\|\tilde{\mathbf{W}}_t - \mathbf{W}_0\|_F + \sqrt{d}) \leq C \Gamma \varepsilon \left(\Gamma \sqrt{\frac{K \log K}{\lambda(\mathbf{C})}} + \sqrt{d}\right). \quad (17)$$

Combining the two bounds above we get

$$\begin{aligned} |f(\mathbf{W}_t, \mathbf{x}) - f(\tilde{\mathbf{W}}_t, \mathbf{c})| &\leq \varepsilon \mathcal{O}\left(\frac{\psi' \Gamma^5 K^2}{\lambda(\mathbf{C})^2} \log\left(\frac{\Gamma \sqrt{n \log K}}{\rho}\right)^3 + \Gamma \left(\Gamma \sqrt{\frac{K \log K}{\lambda(\mathbf{C})}} + \sqrt{d}\right)\right) \\ &\leq \varepsilon \mathcal{O}\left(\frac{\psi' \Gamma^5 K^2}{\lambda(\mathbf{C})^2} \log\left(\frac{\Gamma \sqrt{n \log K}}{\rho}\right)^3\right). \end{aligned}$$

Hence, if  $\varepsilon \leq c' \delta \min\left(\frac{\lambda(\mathbf{C})^2}{\psi' \Gamma^5 K^2 \log\left(\frac{\Gamma \sqrt{n \log K}}{\rho}\right)^3}, \frac{1}{\Gamma \sqrt{d}}\right)$ , we obtain that, for all  $\mathbf{x}$ , the associated cluster  $\mathbf{c}$  and true label assigned to cluster  $\mathbf{y}^* = \mathbf{y}^*(\mathbf{c})$ , we have that

$$|f(\mathbf{W}_t, \mathbf{x}) - \mathbf{y}^*| < |f(\tilde{\mathbf{W}}_t, \mathbf{c}) - f(\mathbf{W}_t, \mathbf{x})| + |f(\tilde{\mathbf{W}}_t, \mathbf{c}) - \mathbf{y}^*| \leq 4\rho + \frac{\delta}{8}.$$

Meanwhile, we can upper bound

$$|\bar{\mathbf{y}}_x^t - \mathbf{y}_x^*| \leq (1 - \alpha_t)|\mathbf{y}_x - \mathbf{y}_x^*| + \alpha_t |f(\mathbf{W}_t, \mathbf{x}) - \mathbf{y}^*| \leq (1 - \alpha_t)|\mathbf{y}_x - \mathbf{y}_x^*| + \alpha_t(4\rho + \frac{\delta}{8}).$$

where  $\bar{\mathbf{y}}_x^t = (1 - \alpha_t)\mathbf{y}_x + \alpha_t f(\mathbf{W}_t, \mathbf{x})$  and  $\mathbf{y}_x^*$  respectively denote the estimated label by our label refinery and the ground truth label of sample  $\mathbf{x}$ . Since  $|\mathbf{y}_x - \mathbf{y}_x^*| < 1$ , by setting  $1 \geq \alpha_t \geq 1 - \frac{3}{4}\delta$  and  $\rho \leq \delta/32$ , we have

$$|\bar{\mathbf{y}}_x^t - \mathbf{y}_x^*| < \frac{\delta}{2}$$

This means that for any sample  $\mathbf{x}_i$ , we have  $|\bar{\mathbf{y}}_i^t - \mathbf{y}_i^*| < \delta/2$ . Therefore, our label refinery gives the correct estimated labels for all samples. By using the same setting, we obtain

$$|f(\mathbf{W}_t, \mathbf{x}) - \mathbf{y}^*| < \delta/2.$$

This means that for any sample  $\mathbf{x}_i$ , we have  $|f(\mathbf{W}_t, \mathbf{x}_i) - \mathbf{y}_i^*| < \delta/2$ . Therefore,  $\mathbf{W}_t$  gives the correct estimated labels for all samples. This completes all proofs.  $\square$

## C.2.4 Proofs of Auxiliary Theories in Appendix C.2

### C.2.5 Proof of Theorem 6

*Proof.* The proof will be done inductively over the properties of gradient descent iterates and is inspired from the recent work [21, 19]. The main difference is that this work uses the label estimation  $\bar{\mathbf{y}}^t = (1 - \alpha_t)\mathbf{y} + \alpha_t f(\mathbf{w}_t)$  and minimizes the squared loss, while both [21, 19] use the corrupted label  $\mathbf{y}$  and then minimize the squared loss. By comparison, our method is much more complicated and gives different proofs. Let us introduce the notation related to the residual. Set  $\mathbf{r}_t = f(\mathbf{w}_t) - \bar{\mathbf{y}}^t$  and let  $\mathbf{r}_0 = f(\mathbf{w}_0) - \bar{\mathbf{y}}^0$  be the initial residual. We keep track of the growth of the residual by partitioning the residual as  $\mathbf{r}_t = \hat{\mathbf{r}}_t + \hat{\mathbf{e}}_t$  where

$$\hat{\mathbf{e}}_t = \mathcal{P}_{S_-}(\mathbf{r}_t) \quad , \quad \hat{\mathbf{r}}_t = \mathcal{P}_{S_+}(\mathbf{r}_t).$$

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

$$\|\hat{\mathbf{e}}_t\|_2 \leq \|\hat{\mathbf{e}}_0\|_2 + \sqrt{n} \sum_{i=0}^t |\alpha_i - \alpha_{i+1}| \leq \|\hat{\mathbf{e}}_0\|_2 + \frac{\psi_1}{2} \|\mathbf{r}_0\|_2, \quad (18)$$

$$\|\hat{\mathbf{r}}_t\|_2^2 \leq \left(1 - \frac{\eta\alpha^2}{4}\right)^t \|\hat{\mathbf{r}}_0\|_2^2 + 2\sqrt{n} \sum_{t=0}^{t-1} \left(1 - \frac{\eta\alpha^2}{4}\right)^{t-t} |\alpha_t - \alpha_{t+1}|, \quad (19)$$

$$\begin{aligned} \frac{\alpha}{4} \|\mathbf{w}_t - \mathbf{w}_0\|_2 + \|\hat{\mathbf{r}}_t\|_2 &\leq \|\hat{\mathbf{r}}_0\|_2 + 2\sqrt{n} \sum_{i=0}^t |\alpha_i - \alpha_{i+1}| \leq \|\mathbf{r}_0\|_2 + 2\sqrt{n} \sum_{i=0}^t |\alpha_i - \alpha_{i+1}| \\ &\leq (1 + \phi) \|\mathbf{r}_0\|_2, \end{aligned} \quad (20)$$

where the last line uses the assumption that  $2\sqrt{n} \lim_{t \rightarrow +\infty} \sum_{i=0}^t |\alpha_i - \alpha_{i+1}| \leq \psi_1 \|\mathbf{r}_0\|_2$ . Assuming these conditions hold till some  $t > 0$ , inductively, we focus on iteration  $t + 1$ . First, note that these conditions imply that for all  $t \geq i \geq 0$ ,  $\mathbf{w}_i \in \mathcal{D}$  where  $\mathcal{D} = \{\mathbf{w} \in \mathbb{R}^P \mid \|\mathbf{w} - \mathbf{w}_0\|_2 \leq \frac{4(1+\psi_1)\|\mathbf{r}_0\|_2}{\alpha}\}$  is the Euclidian ball around  $\mathbf{w}_0$  of radius  $\frac{4(1+\psi_1)\|\mathbf{r}_0\|_2}{\alpha}$ . This directly follows from (20) induction hypothesis. Next, we claim that  $\mathbf{w}_{t+1}$  is still within the set  $\mathcal{D}$ . From Lemma 5, we have that if the results in Eqn. (20) holds, then it holds that

$$\mathbf{w}_{t+1} \in \mathcal{D} = \left\{ \mathbf{w} \in \mathbb{R}^P \mid \|\mathbf{w} - \mathbf{w}_0\|_2 \leq \frac{4(1 + \psi_1)\|\mathbf{r}_0\|_2}{\alpha} \right\}.$$

In this way, we can directly use the results in previous lemmas and assumptions. Then we will prove that (19) and (20) hold for  $t + 1$  as well. Note that, following Lemma 3, gradient descent iterate can be written as

$$\mathbf{r}_{t+1} = (\mathbf{I} - \eta\mathbf{G}(\mathbf{w}_t))\mathbf{r}_t + \bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1}.$$

Since both column and row space of  $\mathbf{G}(\mathbf{w}_t)$  is subset of  $\mathcal{S}_+$ , we have that

$$\hat{\mathbf{e}}_{t+1} = \mathcal{P}_{\mathcal{S}_-}((\mathbf{I} - \eta\mathbf{G}(\mathbf{w}_t))\mathbf{r}_t + \bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1}) \quad (21)$$

$$= \mathcal{P}_{\mathcal{S}_-}(\mathbf{r}_t) + \mathcal{P}_{\mathcal{S}_-}(\bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1}) \quad (22)$$

$$= \hat{\mathbf{e}}_t + \mathcal{P}_{\mathcal{S}_-}(\bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1}) \quad (23)$$

$$= \hat{\mathbf{e}}_t + \mathcal{P}_{\mathcal{S}_-}((\alpha_{t+1} - \alpha_t)\mathbf{y}) \quad (24)$$

$$= \hat{\mathbf{e}}_0 + \sum_{t=0}^t \mathcal{P}_{\mathcal{S}_-}((\alpha_{t+1} - \alpha_t)\mathbf{y}) \quad (25)$$

$$(26)$$

So we can upper bound

$$\|\hat{\mathbf{e}}_t\|_2 \leq \|\hat{\mathbf{e}}_0\|_2 + 2\sqrt{n} \sum_{i=0}^t |\alpha_i - \alpha_{i+1}| \leq \|\hat{\mathbf{e}}_0\|_2 + \psi_1 \|\mathbf{r}_0\|_2. \quad (27)$$

This shows the first statement of the induction. Next, over  $\mathcal{S}_+$ , we have

$$\hat{\mathbf{r}}_{t+1} = \mathcal{P}_{\mathcal{S}_+}((\mathbf{I} - \eta\mathbf{G}(\mathbf{w}_t))\mathbf{r}_t + \bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1}) \quad (28)$$

$$= \mathcal{P}_{\mathcal{S}_+}((\mathbf{I} - \eta\mathbf{G}(\mathbf{w}_t))\hat{\mathbf{r}}_t) + \mathcal{P}_{\mathcal{S}_+}((\mathbf{I} - \eta\mathbf{G}(\mathbf{w}_t))\hat{\mathbf{e}}_t) + \mathcal{P}_{\mathcal{S}_+}(\bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1}) \quad (29)$$

$$= \mathcal{P}_{\mathcal{S}_+}((\mathbf{I} - \eta\mathbf{G}(\mathbf{w}_t))\hat{\mathbf{r}}_t) + \mathcal{P}_{\mathcal{S}_+}(\bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1}) \quad (30)$$

$$= (\mathbf{I} - \eta\mathbf{G}(\mathbf{w}_t))\hat{\mathbf{r}}_t + \bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1} \quad (31)$$

where the second line uses the fact that  $\hat{\mathbf{e}}_t \in \mathcal{S}_-$  and last line uses the fact that  $\hat{\mathbf{r}}_t \in \mathcal{S}_+$ , in the last line, we let  $\hat{\mathbf{y}}_t = \mathcal{P}_{\mathcal{S}_+}(\bar{\mathbf{y}}^t)$ . Then we can rewrite  $\bar{\mathbf{y}}^t - \bar{\mathbf{y}}^{t+1}$  as

$$\begin{aligned} \hat{\mathbf{y}}_t - \hat{\mathbf{y}}_{t+1} &= (1 - \alpha_t)\mathbf{y} + \alpha_t f(\mathbf{w}_t) - (1 - \alpha_{t+1})\mathbf{y} + \alpha_{t+1} f(\mathbf{w}_{t+1}) \\ &= (\alpha_{t+1} - \alpha_t)\mathbf{y} + \alpha_t (f(\mathbf{w}_t) - f(\mathbf{w}_{t+1})) - (\alpha_{t+1} - \alpha_t) f(\mathbf{w}_{t+1}). \end{aligned}$$

At the same time, we can upper bound

$$\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_F = \eta \|\mathcal{J}(\mathbf{w}_t)^T \mathbf{r}_t\|_2 \stackrel{\textcircled{1}}{\leq} \eta \|\mathcal{J}(\mathbf{w}_t)^T \hat{\mathbf{r}}_t\|_2 \leq \eta\beta \|\hat{\mathbf{r}}_t\|_2.$$

In this way, we can obtain

$$\begin{aligned} & \|\hat{\mathbf{r}}_{t+1}\|_2 \\ & \leq \|(\mathbf{I} - \eta\mathbf{G}(\mathbf{w}_t))\hat{\mathbf{r}}_t\|_2 + \|(\alpha_t - \alpha_{t+1})\mathbf{y}\|_2 + \alpha_t \|f(\mathbf{w}_t) - f(\mathbf{w}_{t+1})\|_2 + \|(\alpha_{t+1} - \alpha_t)f(\mathbf{w}_{t+1})\|_2 \\ & \stackrel{\textcircled{1}}{\leq} \left(1 - \frac{\eta\alpha^2}{2}\right) \|\hat{\mathbf{r}}_t\|_2 + \alpha_t \beta \|\mathbf{w}_t - \mathbf{w}_{t+1}\|_2 + 2\sqrt{n} \cdot |\alpha_t - \alpha_{t+1}| \\ & \leq \left(1 - \frac{\eta\alpha^2}{2}\right) \|\hat{\mathbf{r}}_t\|_2 + \alpha_t \beta^2 \eta \|\hat{\mathbf{r}}_t\|_2 + 2\sqrt{n} \cdot |\alpha_t - \alpha_{t+1}| \\ & \stackrel{\textcircled{2}}{\leq} \left(1 - \frac{\eta\alpha^2}{4}\right) \|\hat{\mathbf{r}}_t\|_2 + 2\sqrt{n} \cdot |\alpha_t - \alpha_{t+1}| \end{aligned}$$

where  $\textcircled{1}$  uses in Lemma 6,  $\|\mathbf{y}\|_2 \leq \sqrt{n}$  and  $\|f(\mathbf{w}_{t+1})\|_2 \leq \sqrt{n}$ ,  $\textcircled{2}$  uses  $\alpha_t \leq \frac{\alpha^2}{4\beta^2}$ . This result further yields

$$\|\hat{\mathbf{r}}_t\|_2 \leq \left(1 - \frac{\eta\alpha^2}{4}\right)^t \|\hat{\mathbf{r}}_0\|_2 + 2\sqrt{n} \sum_{t=0}^{t-1} \left(1 - \frac{\eta\alpha^2}{4}\right)^{t-t} |\alpha_t - \alpha_{t+1}|$$

On the other hand, we have

$$\begin{aligned} \|(\mathbf{I} - \eta\mathbf{G}(\mathbf{w}_t))\hat{\mathbf{r}}_t\|_2^2 & \leq \|\hat{\mathbf{r}}_t\|_2^2 - 2\eta\hat{\mathbf{r}}_t^T \mathbf{G}(\mathbf{w}_t)\hat{\mathbf{r}}_t + \eta^2\hat{\mathbf{r}}_t^T \mathbf{G}^T(\mathbf{w}_t)\mathbf{G}(\mathbf{w}_t)\hat{\mathbf{r}}_t \\ & \leq \|\hat{\mathbf{r}}_t\|_2^2 - 2\eta\hat{\mathbf{r}}_t^T \mathcal{J}(\mathbf{w}_t)\mathcal{J}^T(\mathbf{w}_t)\hat{\mathbf{r}}_t + \eta^2\beta^2\hat{\mathbf{r}}_t^T \mathcal{J}(\mathbf{w}_t)\mathcal{J}^T(\mathbf{w}_t)\hat{\mathbf{r}}_t \\ & = \|\hat{\mathbf{r}}_t\|_2^2 - \eta(2 - \eta\beta^2)\|\mathcal{J}^T(\mathbf{w}_t)\hat{\mathbf{r}}_t\|_2^2 \\ & \leq \|\hat{\mathbf{r}}_t\|_2^2 - \eta\|\mathcal{J}^T(\mathbf{w}_t)\hat{\mathbf{r}}_t\|_2^2, \end{aligned}$$

where the last line use  $\eta \leq \frac{1}{\beta^2}$ . This further gives

$$\|(\mathbf{I} - \eta \mathbf{G}(\mathbf{w}_t))\widehat{\mathbf{r}}_t\|_2 \leq \sqrt{\|\widehat{\mathbf{r}}_t\|_2^2 - \eta \|\mathcal{J}^T(\mathbf{w}_t)\widehat{\mathbf{r}}_t\|_2^2} \leq \|\widehat{\mathbf{r}}_t\|_2 - \frac{\eta \|\mathcal{J}^T(\mathbf{w}_t)\widehat{\mathbf{r}}_t\|_2^2}{2\|\widehat{\mathbf{r}}_t\|_2}.$$

Therefore, we can upper bound  $\|\widehat{\mathbf{r}}_t\|_2$  in another way which can help to bound  $\|\mathbf{w}_{t+1} - \mathbf{w}_0\|_2$ :

$$\begin{aligned} & \|\widehat{\mathbf{r}}_{t+1}\|_2 \\ & \leq \|(\mathbf{I} - \eta \mathbf{G}(\mathbf{w}_t))\widehat{\mathbf{r}}_t\|_2 + \|(\alpha_t - \alpha_{t+1})\mathbf{y}\|_2 + (1 - \alpha_t)\|f(\mathbf{w}_t) - f(\mathbf{w}_{t+1})\|_2 + \|(\alpha_{t+1} - \alpha_t)f(\mathbf{w}_{t+1})\|_2 \\ & \leq \|(\mathbf{I} - \eta \mathbf{G}(\mathbf{w}_t))\widehat{\mathbf{r}}_t\|_2 + (1 - \alpha_t)\beta\|\mathbf{w}_t - \mathbf{w}_{t+1}\|_2 + 2\sqrt{n} \cdot |\alpha_t - \alpha_{t+1}| \\ & = \|(\mathbf{I} - \eta \mathbf{G}(\mathbf{w}_t))\widehat{\mathbf{r}}_t\|_2 + (1 - \alpha_t)\beta\eta\|\mathcal{J}^T(\mathbf{w}_t)\mathbf{r}_t\|_2 + 2\sqrt{n} \cdot |\alpha_t - \alpha_{t+1}| \\ & \leq \|\widehat{\mathbf{r}}_t\|_2 - \frac{\eta \|\mathcal{J}^T(\mathbf{w}_t)\widehat{\mathbf{r}}_t\|_2^2}{2\|\widehat{\mathbf{r}}_t\|_2} + (1 - \alpha_t)\beta\eta\|\mathcal{J}^T(\mathbf{w}_t)\mathbf{r}_t\|_2 + 2\sqrt{n} \cdot |\alpha_t - \alpha_{t+1}|. \end{aligned}$$

Since the distance of  $\mathbf{w}_{t+1}$  to initial point satisfies :

$$\|\mathbf{w}_{t+1} - \mathbf{w}_0\|_2 \leq \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2 + \|\mathbf{w}_t - \mathbf{w}_0\|_2 \leq \|\mathbf{w}_t - \mathbf{w}_0\|_2 + \eta\|\mathcal{J}^T(\mathbf{w}_t)\mathbf{r}_t\|_2,$$

we can further bound

$$\begin{aligned} & \frac{\alpha}{4}\|\mathbf{w}_{t+1} - \mathbf{w}_0\|_2 + \|\widehat{\mathbf{r}}_{t+1}\|_2 \\ & \leq \frac{\alpha}{4}(\|\mathbf{w}_t - \mathbf{w}_0\|_2 + \eta\|\mathcal{J}^T(\mathbf{w}_t)\mathbf{r}_t\|_2) + \|\widehat{\mathbf{r}}_t\|_2 - \frac{\eta \|\mathcal{J}^T(\mathbf{w}_t)\widehat{\mathbf{r}}_t\|_2^2}{2\|\widehat{\mathbf{r}}_t\|_2} \\ & \quad + (1 - \alpha_t)\beta\eta\|\mathcal{J}^T(\mathbf{w}_t)\mathbf{r}_t\|_2 + 2\sqrt{n} \cdot |\alpha_t - \alpha_{t+1}| \\ & \leq \frac{\alpha}{4}\|\mathbf{w}_t - \mathbf{w}_0\|_2 + \|\widehat{\mathbf{r}}_t\|_2 + \frac{\eta}{4}\|\mathcal{J}^T(\mathbf{w}_t)\mathbf{r}_t\|_2 \left( \alpha + 4(1 - \alpha_t)\beta - 2\frac{\|\mathcal{J}^T(\mathbf{w}_t)\widehat{\mathbf{r}}_t\|_2}{\|\widehat{\mathbf{r}}_t\|_2} \right) + 2\sqrt{n} \cdot |\alpha_t - \alpha_{t+1}| \\ & \stackrel{\textcircled{1}}{\leq} \frac{\alpha}{4}\|\mathbf{w}_t - \mathbf{w}_0\|_2 + \|\widehat{\mathbf{r}}_t\|_2 + 2\sqrt{n} \cdot |\alpha_t - \alpha_{t+1}| \\ & \leq \|\widehat{\mathbf{r}}_0\|_2 + 2\sqrt{n} \sum_{i=0}^t |\alpha_i - \alpha_{i+1}| \leq \|\mathbf{r}_0\|_2 + 2\sqrt{n} \sum_{i=0}^t |\alpha_i - \alpha_{i+1}|, \end{aligned}$$

where  $\textcircled{1}$  uses  $\frac{\|\mathcal{J}^T(\mathbf{w}_t)\widehat{\mathbf{r}}_t\|_2}{\|\widehat{\mathbf{r}}_t\|_2} \geq \alpha$  and  $\alpha_t \leq \frac{\alpha}{4\beta}$ .

By setting  $t \geq \frac{5}{\eta\alpha^2} \log\left(\frac{\|\mathbf{r}_0\|_2}{(1 - \alpha_{\max})^\nu}\right)$  and  $\frac{\eta\alpha^2}{4} \leq \frac{\eta\beta^2}{4} \leq \frac{1}{8}$  where  $\alpha_{\max} = \max_t \alpha_t$ , then we have  $\log \frac{1}{1 - \frac{\eta\alpha^2}{4}} \geq \log\left(1 + \frac{\eta\alpha^2}{4}\right) \geq \frac{\eta\alpha^2}{5}$  and thus

$$\left(1 - \frac{\eta\alpha^2}{4}\right)^t \|\widehat{\mathbf{r}}_0\|_2 \leq \left(1 - \frac{\eta\alpha^2}{4}\right)^t \|\mathbf{r}_0\|_2 \leq \left(1 - \frac{\eta\alpha^2}{4}\right)^t \|\mathbf{r}_0\|_2 \leq (1 - \alpha_{\max})^\nu.$$

In this way, we can further obtain

$$\|\widehat{\mathbf{r}}_t\|_\infty \leq \|\widehat{\mathbf{r}}_t\|_2 \leq (1 - \alpha_{\max})^\nu + 2\sqrt{n} \sum_{i=0}^{t-1} \left(1 - \frac{\eta\alpha^2}{4}\right)^{t-i} |\alpha_i - \alpha_{i+1}|$$

and

$$\begin{aligned} (1 - \alpha_t)\|\mathcal{P}_{S_+}(f(\mathbf{w}_t) - \mathbf{y})\|_\infty & = \|\mathcal{P}_{S_+}(f(\mathbf{w}_t) - (1 - \alpha_t)\mathbf{y} - \alpha_t f(\mathbf{w}_t))\|_\infty = \|\widehat{\mathbf{r}}_t\|_\infty \leq \|\widehat{\mathbf{r}}_t\|_2 \\ & \leq (1 - \alpha_{\max})^\nu + 2\sqrt{n} \sum_{t=0}^{t-1} \left(1 - \frac{\eta\alpha^2}{4}\right)^{t-t} |\alpha_t - \alpha_{t+1}| \end{aligned}$$

Finally, we can obtain the desired results:

$$\begin{aligned} \|f(\mathbf{w}_t) - \mathbf{y}^*\|_\infty & \stackrel{\textcircled{1}}{=} \|\mathcal{P}_{S_+}(f(\mathbf{w}_t)) - \mathcal{P}_{S_+}(\mathbf{y}^*)\|_\infty \\ & \leq \|\mathcal{P}_{S_+}(f(\mathbf{w}_t) - \mathbf{y})\|_\infty + \|\mathcal{P}_{S_+}(\mathbf{y} - \mathbf{y}^*)\|_\infty \\ & \leq 2\nu + \frac{2\sqrt{n}}{1 - \alpha_t} \sum_{i=0}^{t-1} \left(1 - \frac{\eta\alpha^2}{4}\right)^{t-i} |\alpha_i - \alpha_{i+1}|, \end{aligned}$$

where ① holds since  $f(\mathbf{w}_t) - \mathbf{y}^* \in \mathcal{S}_+$  and  $\|\mathcal{P}_{\mathcal{S}_+}(f(\mathbf{w}_t) - \mathbf{y})\|_\infty = \|\mathcal{P}_{\mathcal{S}_+}(f(\mathbf{w}_t) - \mathbf{y})\|_\infty$ . If  $e$  is  $s$  sparse and  $\mathcal{S}_+$  is diffused, applying Definition 3 we have

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

The proof is completed.  $\square$

### C.2.6 Proof of Theorem 7

*Proof.* The proof is based on the meta Theorem 6, hence we need to verify its Assumptions 2 and 3 with proper values and apply Lemma 8 to get  $\|\mathcal{P}_{\mathcal{S}_+}(e)\|_\infty$ . We will also make significant use of Corollary 4.

Using Corollary 4, Assumption 3 holds with  $L = \Gamma \sqrt{\frac{c_{up}n}{kK}} \|\mathbf{C}\|$  where  $L$  is the Lipschitz constant of Jacobian spectrum. Denote Using Lemma 7 with probability  $1 - K^{-100}$ , we have that  $\|\mathbf{r}_0\|_2 = \|\bar{\mathbf{y}}^0 - f(\mathbf{W}_0)\|_2 = \|\mathbf{y} - f(\mathbf{W}_0)\|_2 \leq \Gamma \sqrt{c_0 n \log K / 128}$  for some  $c_0 > 0$ . Corollary 4 guarantees a uniform bound for  $\beta$ , hence in Assumption 2, we pick

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

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

$$\alpha = \frac{\alpha'}{2} \quad \text{where} \quad \alpha' = \sqrt{\frac{c_{low}n\lambda(\mathbf{C})}{2K}},$$

We wish to verify Assumption 2 over the radius of

$$R = \frac{4\|f(\mathbf{W}_0) - \mathbf{y}\|_2}{\alpha} \leq \frac{\Gamma \sqrt{c_0 n \log K / 8}}{\alpha} = \Gamma \sqrt{\frac{c_0 n \log K / 2}{\frac{c_{low}n\lambda(\mathbf{C})}{2K}}} = \Gamma \sqrt{\frac{c_0 K \log K}{c_{low}\lambda(\mathbf{C})}},$$

neighborhood of  $\mathbf{W}_0$ . What remains is ensuring that Jacobian over  $\mathcal{S}_+$  is lower bounded by  $\alpha$ . Our choice of  $k$  guarantees that at the initialization, with probability  $1 - K^{-100}$ , we have

$$\sigma(\mathcal{J}(\mathbf{W}_0, \mathbf{X}), \mathcal{S}_+) \geq \alpha'.$$

Suppose  $LR \leq \alpha = \alpha'/2$  which can be achieved by using large  $k$ . Using triangle inequality on Jacobian spectrum, for any  $\mathbf{W} \in \mathcal{D}$ , using  $\|\mathbf{W} - \mathbf{W}_0\|_F \leq R$ , we would have

$$\sigma(\mathcal{J}(\mathbf{W}, \mathbf{X}), \mathcal{S}_+) \geq \sigma(\mathcal{J}(\mathbf{W}_0, \mathbf{X}), \mathcal{S}_+) - LR \geq \alpha' - \alpha = \alpha.$$

Now, observe that

$$LR = (1 + \psi_1) \Gamma \sqrt{\frac{c_{up}n}{kK}} \|\mathbf{C}\| \Gamma \sqrt{\frac{c_0 K \log(K)}{c_{low}\lambda(\mathbf{C})}} = (1 + \psi_1) \Gamma^2 \|\mathbf{C}\| \sqrt{\frac{c_{up}c_0 n \log K}{c_{low}k\lambda(\mathbf{C})}} \quad (32)$$

$$\leq \frac{\alpha'}{2} = \sqrt{\frac{c_{low}n\lambda(\mathbf{C})}{8K}}, \quad (33)$$

as  $k$  satisfies

$$k \geq \mathcal{O} \left( (1 + \psi_1)^2 \Gamma^4 \|\mathbf{C}\|^2 \frac{c_{up}K \log(K)}{c_{low}^2 \lambda(\mathbf{C})^2} \right) \geq \mathcal{O} \left( \frac{(1 + \psi_1)^2 \Gamma^4 K \log(K) \|\mathbf{C}\|^2}{\lambda(\mathbf{C})^2} \right).$$

Finally, since  $LR = 4(1 + \psi_1)L\|\mathbf{r}_0\|_2/\alpha \leq \alpha$ , the learning rate is

$$\eta \leq \frac{1}{2\beta^2} \min(1, \frac{\alpha\beta}{L\|\mathbf{r}_0\|_2}) = \frac{1}{2\beta^2} = \frac{K}{2c_{up}n\Gamma^2 \|\mathbf{C}\|^2}.$$

Overall, the assumptions of Theorem 6 holds with stated  $\alpha, \beta, L$  with probability  $1 - 2K^{-100}$  (union bounding initial residual and minimum singular value events). This implies for all  $t > 0$  the distance of current iterate to initial obeys

$$\|\mathbf{W}_t - \mathbf{W}_0\|_F \leq R.$$

The final step is the properties of the label corruption. Using Lemma 8, we find that

$$\|\mathcal{P}_{S_+}(\mathbf{y}^* - \mathbf{y})\|_\infty \leq 2\rho.$$

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

$$\frac{5}{\eta\alpha^2} \log\left(\frac{\|\mathbf{r}_0\|_2}{2(1 - \alpha_{\max})\rho}\right) \leq \frac{5}{\eta\alpha^2} \log\left(\frac{\Gamma\sqrt{c_0 n \log K/32}}{2(1 - \alpha_{\max})\rho}\right) = \mathcal{O}\left(\frac{K}{\eta n \lambda(\mathbf{C})} \log\left(\frac{\Gamma\sqrt{n \log K}}{(1 - \alpha_{\max})\rho}\right)\right) \leq t,$$

denoting the clean labels by  $\tilde{\mathbf{y}}$  and applying Theorem 6, we have that, the infinity norm of the residual obeys (using  $\|\mathcal{P}_{S_+}(\mathbf{e})\|_\infty = \|\mathcal{P}_{S_+}(\mathbf{y} - \mathbf{y}^*)\|_\infty \leq 2\rho$ )

$$\|f(\mathbf{W}) - \mathbf{y}^*\|_\infty \leq 4\rho.$$

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

### C.2.7 Proof of Theorem 8

*Proof.* Since  $\tilde{\mathbf{W}}_t$  are the noiseless iterations, with probability  $1 - 2K^{-100}$ , the statements of Theorem 7 hold on  $\tilde{\mathbf{W}}_t$ . To proceed with proof, we first introduce short hand notations. We use

$$\mathbf{r}_i = f(\mathbf{W}_i, \mathbf{X}) - \tilde{\mathbf{y}}^i, \tilde{\mathbf{r}}_i = f(\tilde{\mathbf{W}}_i, \tilde{\mathbf{X}}_i) - \tilde{\mathbf{y}}^i \quad (34)$$

$$\mathcal{J}_i = \mathcal{J}(\mathbf{W}_i, \mathbf{X}), \mathcal{J}_{i+1,i} = \mathcal{J}(\mathbf{W}_{i+1}, \mathbf{W}_i, \mathbf{X}), \tilde{\mathcal{J}}_i = \mathcal{J}(\tilde{\mathbf{W}}_i, \tilde{\mathbf{X}}), \tilde{\mathcal{J}}_{i+1,i} = \mathcal{J}(\tilde{\mathbf{W}}_{i+1}, \tilde{\mathbf{W}}_i, \tilde{\mathbf{X}}) \quad (35)$$

$$d_i = \|\mathbf{W}_i - \tilde{\mathbf{W}}_i\|_F, p_i = \|\mathbf{r}_i - \tilde{\mathbf{r}}_i\|_2, \beta = \Gamma\|\mathbf{C}\|\sqrt{c_{up}n/K}, L = \Gamma\|\mathbf{C}\|\sqrt{c_{up}n/Kk}. \quad (36)$$

Here  $\beta$  is the upper bound on the Jacobian spectrum and  $L$  is the spectral norm Lipschitz constant as in Theorem 4. Applying Lemma 9, note that

$$\|\mathcal{J}(\mathbf{W}_t, \mathbf{X}) - \mathcal{J}(\tilde{\mathbf{W}}_t, \tilde{\mathbf{X}})\| \leq L\|\tilde{\mathbf{W}}_t - \mathbf{W}_t\|_2 + \Gamma\sqrt{n}\varepsilon \leq Ld_t + \Gamma\sqrt{n}\varepsilon \quad (37)$$

$$\|\mathcal{J}(\mathbf{W}_{t+1}, \mathbf{W}_t, \mathbf{X}) - \mathcal{J}(\tilde{\mathbf{W}}_{t+1}, \tilde{\mathbf{W}}_t, \tilde{\mathbf{X}})\| \leq L(d_t + d_{t+1})/2 + \Gamma\sqrt{n}\varepsilon. \quad (38)$$

By defining

$$\hat{\mathbf{e}}_t = \mathcal{P}_{S_-}(\tilde{\mathbf{r}}_t) \quad , \quad \hat{\mathbf{r}}_t = \mathcal{P}_{S_+}(\tilde{\mathbf{r}}_t),$$

then we can use Theorem 6 and the assumption that  $2\sqrt{n}\sum_{i=0}^{t-1}\left(1 - \frac{\eta\alpha^2}{4}\right)^{t-i}|\alpha_i - \alpha_{i+1}| \leq \psi_2\|\mathbf{r}_0\|_2^2$  to obtain

$$\|\hat{\mathbf{e}}_t\|_2 \leq \|\hat{\mathbf{e}}_0\|_2 + \sqrt{n}\sum_{i=0}^t|\alpha_i - \alpha_{i+1}| \leq \|\hat{\mathbf{e}}_0\|_2 + \frac{\psi}{2}\|\mathbf{r}_0\|_2, \quad (39)$$

$$\|\hat{\mathbf{r}}_t\|_2^2 \leq \left(1 - \frac{\eta\alpha^2}{4}\right)^t \|\hat{\mathbf{r}}_0\|_2^2 + 2\sqrt{n}\sum_{i=0}^{t-1}\left(1 - \frac{\eta\alpha^2}{4}\right)^{t-i}|\alpha_i - \alpha_{i+1}| \leq \|\hat{\mathbf{r}}_0\|_2^2 + \psi_2\|\mathbf{r}_0\|_2^2. \quad (40)$$

Therefore, we can upper bound

$$\|\tilde{\mathbf{r}}_t\|_2 = \|\hat{\mathbf{e}}_t\|_2 + \|\hat{\mathbf{r}}_t\|_2 \leq \|\hat{\mathbf{e}}_0\|_2 + \frac{\psi}{2}\|\mathbf{r}_0\|_2 + \|\hat{\mathbf{r}}_0\|_2 + \sqrt{\psi_2}\|\mathbf{r}_0\|_2 = \left(1 + \frac{\psi}{2} + \sqrt{\psi_2}\right)\|\mathbf{r}_0\|_2. \quad (41)$$

Following this and setting  $\|\tilde{\mathbf{r}}_t\|_2 \leq \psi'\|\mathbf{r}_0\|_2$ , note that parameter satisfies

$$\mathbf{W}_{i+1} = \mathbf{W}_i - \eta\mathcal{J}_i\mathbf{r}_i \quad , \quad \tilde{\mathbf{W}}_{i+1} = \tilde{\mathbf{W}}_i - \eta\tilde{\mathcal{J}}_i^T\tilde{\mathbf{r}}_i \quad (42)$$

$$\|\mathbf{W}_{i+1} - \tilde{\mathbf{W}}_{i+1}\|_F \leq \|\mathbf{W}_i - \tilde{\mathbf{W}}_i\|_F + \eta\|\mathcal{J}_i - \tilde{\mathcal{J}}_i\|\|\mathbf{r}_i\|_F + \eta\|\mathcal{J}_i\|\|\mathbf{r}_i - \tilde{\mathbf{r}}_i\|_2 \quad (43)$$

$$d_{i+1} \leq d_i + \eta(\psi'(Ld_i + \Gamma\sqrt{n}\varepsilon)\|\mathbf{r}_0\|_2 + \beta p_i), \quad (44)$$

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

$$\mathbf{r}_{i+1} = \mathbf{r}_i - \eta \mathcal{J}_{i+1,i} \mathcal{J}_i^T \mathbf{r}_i \implies \quad (45)$$

$$\mathbf{r}_{i+1} - \tilde{\mathbf{r}}_{i+1} \quad (46)$$

$$= (\mathbf{r}_i - \tilde{\mathbf{r}}_i) - \eta(\mathcal{J}_{i+1,i} - \tilde{\mathcal{J}}_{i+1,i}) \mathcal{J}_i^T \mathbf{r}_i - \eta \tilde{\mathcal{J}}_{i+1,i} (\mathcal{J}_i^T - \tilde{\mathcal{J}}_i^T) \mathbf{r}_i - \eta \tilde{\mathcal{J}}_{i+1,i} \tilde{\mathcal{J}}_i^T (\mathbf{r}_i - \tilde{\mathbf{r}}_i). \quad (47)$$

$$\mathbf{r}_{i+1} - \tilde{\mathbf{r}}_{i+1} = (\mathbf{I} - \eta \tilde{\mathcal{J}}_{i+1,i} \tilde{\mathcal{J}}_i^T) (\mathbf{r}_i - \tilde{\mathbf{r}}_i) - \eta(\mathcal{J}_{i+1,i} - \tilde{\mathcal{J}}_{i+1,i}) \mathcal{J}_i^T \mathbf{r}_i - \eta \tilde{\mathcal{J}}_{i+1,i} (\mathcal{J}_i^T - \tilde{\mathcal{J}}_i^T) \mathbf{r}_i. \quad (48)$$

$$\|\mathbf{r}_{i+1} - \tilde{\mathbf{r}}_{i+1}\|_2 \leq \|\mathbf{r}_i - \tilde{\mathbf{r}}_i\|_2 + \eta\beta \|\mathbf{r}_i\|_2 (L(3d_t + d_{t+1})/2 + 2\Gamma\sqrt{n}\varepsilon). \quad (49)$$

$$\|\mathbf{r}_{i+1} - \tilde{\mathbf{r}}_{i+1}\|_2 \leq \|\mathbf{r}_i - \tilde{\mathbf{r}}_i\|_2 + \eta\beta (\|\tilde{\mathbf{r}}_0\|_2 + p_i) (L(3d_t + d_{t+1})/2 + 2\Gamma\sqrt{n}\varepsilon). \quad (50)$$

where we used  $\|\mathbf{r}_i\|_2 \leq p_i + \psi' \|\mathbf{r}_0\|_2$  and  $\|(\mathbf{I} - \eta \tilde{\mathcal{J}}_{i+1,i} \tilde{\mathcal{J}}_i^T) \mathbf{v}\|_2 \leq \|\mathbf{v}\|_2$  which follows from Lemma 6. This implies

$$p_{i+1} \leq p_i + \eta\beta (\psi' \|\mathbf{r}_0\|_2 + p_i) (L(3d_t + d_{t+1})/2 + 2\Gamma\sqrt{n}\varepsilon). \quad (51)$$

**Finalizing proof:** Next, using Lemma 7, we have  $\|\mathbf{r}_0\|_2 \leq \Theta := C_0 \Gamma \sqrt{n \log K}$ . We claim that if

$$\boxed{\varepsilon \leq \mathcal{O}\left(\frac{1}{t_0 \eta \Gamma^2 n}\right) \leq \frac{1}{8t_0 \eta \beta \Gamma \sqrt{n}} \quad \text{and} \quad L \leq \frac{2}{5t_0 \eta \Theta (1 + 8\eta t_0 \beta^2)} \leq \frac{1}{30(t_0 \eta \beta)^2 \Theta}}, \quad (52)$$

(where we used  $\eta t_0 \beta^2 \geq 1$ ), for all  $t \leq t_0$ , we have that

$$p_t \leq 8t(1 + \psi') \eta \Gamma \sqrt{n} \varepsilon \Theta \beta \leq \Theta, \quad d_t \leq 2t \eta \Gamma \sqrt{n} \varepsilon \Theta (\psi' + 8\eta t_0 \beta^2). \quad (53)$$

The proof is by induction. Suppose it holds until  $t \leq t_0 - 1$ . At  $t + 1$ , via (44) we have that

$$\frac{d_{t+1} - d_t}{\eta} \leq \psi' (L d_t \Theta + \Gamma \sqrt{n} \varepsilon \Theta) + 8t_0 \eta \beta^2 \Gamma \sqrt{n} \varepsilon \Theta \stackrel{?}{\leq} 2\Gamma \sqrt{n} \varepsilon \Theta (\psi' + 8\eta t_0 \beta^2).$$

Right hand side holds since  $L \leq \frac{1}{2\eta t_0 \Theta}$ . This establishes the induction for  $d_{t+1}$ .

Next, we show the induction on  $p_t$ . Observe that  $3d_t + d_{t+1} \leq 10t_0 \eta \Gamma \sqrt{n} \varepsilon \Theta (\psi' + 8\eta t_0 \beta^2)$ . Following (51) and using  $p_t \leq \Theta$ , we need

$$\frac{p_{t+1} - p_t}{\eta} \leq \beta(1 + \psi') \Theta (L(3d_t + d_{t+1}) + 4\Gamma \sqrt{n} \varepsilon) \stackrel{?}{\leq} \frac{8}{\alpha_{\max}} (1 + \psi') \Gamma \sqrt{n} \varepsilon \Theta \beta \iff \quad (54)$$

$$L(3d_t + d_{t+1}) + 4\Gamma \sqrt{n} \varepsilon \stackrel{?}{\leq} \frac{8}{\alpha_{\max}} \Gamma \sqrt{n} \varepsilon \iff \quad (55)$$

$$L(3d_t + d_{t+1}) \stackrel{?}{\leq} \frac{4}{\alpha_{\max}} \Gamma \sqrt{n} \varepsilon \iff \quad (56)$$

$$10\alpha_{\max} L t_0 \eta (1 + 8\eta t_0 \beta^2) \Theta \stackrel{?}{\leq} 4 \iff \quad (57)$$

$$L \stackrel{?}{\leq} \frac{2}{5t_0 \alpha_{\max} \eta (1 + 8\eta t_0 \beta^2) \Theta}, \quad (58)$$

where  $\alpha_{\max} = \max_{1 \leq t \leq t_0} \alpha_t$ . Concluding the induction since  $L$  satisfies the final line. Consequently, for all  $0 \leq t \leq t_0$ , we have that

$$\begin{aligned} p_t &= \|\mathbf{r}_i - \tilde{\mathbf{r}}_i\|_2 = \|f(\mathbf{W}_i, \mathbf{X}) - \bar{\mathbf{y}}^i - f(\tilde{\mathbf{W}}_i, \tilde{\mathbf{X}}_i) + \tilde{\mathbf{y}}^i\|_2 \\ &\stackrel{\textcircled{1}}{\leq} \alpha_{\max} \|f(\mathbf{W}_i, \mathbf{X}) - f(\tilde{\mathbf{W}}_i, \tilde{\mathbf{X}}_i)\|_2 \\ &\leq 8t(1 + \psi') \eta \Gamma \sqrt{n} \varepsilon \Theta \beta = c_0 t (1 + \psi') \eta \varepsilon \Gamma^3 n^{3/2} \sqrt{\log K}. \end{aligned}$$

where  $\textcircled{1}$  uses the definition of  $\bar{\mathbf{y}}^i = (1 - \alpha_i) \mathbf{y} + \alpha_i f(\mathbf{W}_i, \mathbf{X})$  and  $\tilde{\mathbf{y}}^i = (1 - \alpha_i) \mathbf{y} + \alpha_i f(\tilde{\mathbf{W}}_i, \tilde{\mathbf{X}})$ . In this way, we can obtain

$$\|f(\mathbf{W}_i, \mathbf{X}) - f(\tilde{\mathbf{W}}_i, \tilde{\mathbf{X}}_i)\|_2 \leq c_0 t (1 + \psi') \eta \varepsilon \Gamma^3 n^{3/2} \sqrt{\log K}.$$

Next, note that, condition on  $L$  is implied by

$$k \geq 1000\Gamma^2 n(t_0\eta\beta)^4 \Theta^2 / \alpha_{\max}^2 \quad (59)$$

$$= \mathcal{O} \left( \Gamma^4 n \frac{K^4}{\alpha_{\max}^2 n^4 \lambda(\mathcal{C})^4} \log \left( \frac{\Gamma \sqrt{n \log K}}{\rho} \right)^4 (\|\mathcal{C}\| \Gamma \sqrt{n/K})^4 (\Gamma \sqrt{n \log K})^2 \right) \quad (60)$$

$$= \mathcal{O} \left( \Gamma^{10} \frac{K^2 \|\mathcal{C}\|^4}{\alpha_{\max}^2 \lambda(\mathcal{C})^4} \log \left( \frac{\Gamma \sqrt{n \log K}}{\rho} \right)^4 \log^2(K) \right) \quad (61)$$

which is implied by  $k \geq \mathcal{O} \left( \Gamma^{10} \frac{K^2 \|\mathcal{C}\|^4}{\alpha_{\max}^2 \lambda(\mathcal{C})^4} \log \left( \frac{\Gamma \sqrt{n \log K}}{\rho} \right)^6 \right)$ .

Finally, following (53), distance satisfies

$$d_t \leq 20t\psi'\eta^2 t_0 \Gamma \sqrt{n} \varepsilon \Theta \beta^2 \leq \mathcal{O} \left( t\psi'\eta \varepsilon \frac{\Gamma^4 K n}{\lambda(\mathcal{C})} \log \left( \frac{\Gamma \sqrt{n \log K}}{\rho} \right)^2 \right).$$

The proof is completed.  $\square$

## References

- [1] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pages 9729–9738, 2020. [1](#), [2](#), [3](#)
- [2] X. Wang and G. Qi. Contrastive learning with stronger augmentations. 2021. [1](#)
- [3] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *Proc. Int'l Conf. Machine Learning*, 2020. [2](#)
- [4] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pages 770–778, 2016. [2](#)
- [5] I. Loshchilov and F. Hutter. SGDR: Stochastic gradient descent with warm restarts. In *Int'l Conf. Learning Representations*, 2016. [2](#)
- [6] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *Int'l Conf. Learning Representations*, 2015. [2](#)
- [7] K. Lee, Y. Zhu, K. Sohn, C. Li, J. Shin, and H. Lee. i-mix: A strategy for regularizing contrastive representation learning. *arXiv preprint arXiv:2010.08887*, 2020. [2](#)
- [8] M. Everingham, G. Van, C. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *Int'l. J. Computer Vision*, 88(2):303–338, 2010. [3](#)
- [9] T. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. Zitnick. Microsoft coco: Common objects in context. In *Proc. European Conf. Computer Vision*, pages 740–755. Springer, 2014. [3](#)
- [10] Y. Wu, A. Kirillov, F. Massa, W. Lo, and R. Girshick. Detectron2, 2019. [3](#)
- [11] A. Maurer and M. Pontil. Empirical bernstein bounds and sample variance penalization. *arXiv preprint arXiv:0907.3740*, 2009. [3](#)
- [12] S. Liang and R. Srikant. Why deep neural networks for function approximation? *arXiv preprint arXiv:1610.04161*, 2016. [3](#)
- [13] Z. Lu, H. Pu, F. Wang, Z. Hu, and L. Wang. The expressive power of neural networks: A view from the width. In *Proc. Conf. Neural Information Processing Systems*, pages 6231–6239, 2017. [3](#), [5](#)
- [14] M. Telgarsky. Benefits of depth in neural networks. In *Conf. on Learning Theory*, 2016. [3](#), [5](#)
- [15] N. Cohen, O. Sharir, and A. Shashua. On the expressive power of deep learning: A tensor analysis. In *Conf. on Learning Theory*, pages 698–728, 2016. [3](#), [5](#)
- [16] R. Eldan and O. Shamir. The power of depth for feedforward neural networks. In *Conf. on Learning Theory*, pages 907–940, 2016. [3](#), [5](#)
- [17] L. Sagun, U. Evci, V. Guney, Y. Dauphin, and L. Bottou. Empirical analysis of the hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017. [8](#)

- [18] S. Hochreiter and J. Schmidhuber. Flat minima. *Neural Computation*, 9(1):1–42, 1997. 8
- [19] M. Li, M. Soltanolkotabi, and S. Oymak. Gradient descent with early stopping is provably robust to label noise for overparameterized neural networks. In *Proc. Int'l Conf. Artificial Intelligence and Statistics*, pages 4313–4324, 2020. 9, 10, 11, 12, 13, 17
- [20] S. Oymak and M. Soltanolkotabi. Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 2020. 10
- [21] S. Oymak and M. Soltanolkotabi. Overparameterized nonlinear learning: Gradient descent takes the shortest path? In *Proc. Int'l Conf. Machine Learning*, pages 4951–4960, 2019. 13, 17